Contents

	0.1	Preface	i
		0.1.1 Course Concept	i
		0.1.2 Course Contents	i
		0.1.3 Programming Exercises and JuptyterLab as a Web IDE	ii
		0.1.4 This Document	iii
		0.1.5 Acknowledgments	iii
	0.2	Recorded Syllabus	iii
	0.2		.11
1	\mathbf{Pre}	liminaries	1
	1.1	Administrativa	1
	1.2	Getting Most out of IWGS	3
	1.3	Learning Resources for IWGS	5
	1.4	Goals, Culture, & Outline of the Course	7
	1.5	ALeA – AI-Supported Learning	8
2	Intr	roduction to Programming	7
-	2.1	What is Programming?	17
	$\frac{2.1}{2.2}$	Programming in IWGS	20
	$\frac{2.2}{2.3}$	Programming in Python	20
	2.0	2 3 1 Hello IWGS)2
		2.3.2 JupyterLab a Python Web IDE for IWCS	24
		2.3.2 Supportably a Types	24
		2.3.4 Python Control Structures	20
	24	Some Thoughts about Computers and Programs	26
	2.4 2.5	More about Puthon	20
	2.0	25.1 Company and Iteration	0
		2.5.1 Sequences and Iteration)0 41
		2.5.2 Input and Output	±1 49
		2.5.5 Functions and Libraries in Python	15
		2.5.4 A Final word on Programming in IWGS	19
3	Nu	mbers, Characters, and Strings 4	17
	3.1	Representing and Manipulating Numbers 4	17
	3.2	Characters and their Encodings: ASCII and UniCode	<i>i</i> 1
	3.3	More on Computing with Strings	55
	3.4	More on Functions in Python	58
	3.5	Regular Expressions: Patterns in Strings	52
4	Doc	cuments as Digital Objects 6	57
	4.1	Representing & Manipulating Documents on a Computer	37
	4.2	Measuring Sizes of Documents/Units of Information	70
	4.3	Hypertext Markup Language	72
		4.3.1 Introduction	72
		4.3.2 Interacting with HTML in Web Broswers	74
			_

CONTENTS	CON	TENTS
----------	-----	-------

		4.3.3	A Worked Example: The Contact Form	76
	4.4	Docum	ents as Trees	79
	4.5	An Ove	rview over XML Technologies	84
		4.5.1	Introduction to XML	84
		4.5.2	Computing with XML in Python	87
		4.5.3	XML Namespaces	91
		4.5.4	XPath: Specifying XML Subtrees	92
5	We	b Appli	cations	95
	5.1	Web Ap	pplications: The Idea	95
	5.2	Basic C	oncepts of the World Wide Web	96
		5.2.1	Preliminaries	96
		5.2.2	Addressing on the World Wide Web	97
		5.2.3	Running the World Wide Web	100
	5.3	Recap:	HTML Forms Data Transmission	102
	5.4	Generat	ing HTML on the Server	105
		5.4.1	Routing and Argument Passing in Bottle	106
		5.4.2	Templating in Python via STPL	109
		5.4.3	Completing the Contact Form	111
6	From	nt-end 7	Technologies 1	115
6	Fro : 6.1	nt-end 7 Dynam	Icehnologies 1 ic HTML: Client-side Manipulation of HTML Documents	1 15 115
6	Fro : 6.1	nt-end 7 Dynam 6.1.1	Technologies 1 Ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1	115 115 116
6	Fro 6.1 6.2	nt-end 7 Dynam 6.1.1 Cascadi	Iechnologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1	115 115 116 121
6	Fro : 6.1 6.2	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1	Technologies 1 Ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1	115 115 116 121 121
6	Fro : 6.1 6.2	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2	Technologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1	115 116 121 121 121 123
6	Fro : 6.1 6.2	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3	Technologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 A small but useful Fragment of CSS 1	115 116 121 121 123 125
6	Fro : 6.1 6.2	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4	Iechnologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 A small but useful Fragment of CSS 1	115 116 121 121 123 125 130
6	From 6.16.26.3	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4 jQuery:	Technologies 1 Ac HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 Ing Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 A small but useful Fragment of CSS 1 Write Less, Do More 1	115 116 121 121 123 125 130 131
6 7	 From 6.1 6.2 6.3 Pra 	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4 jQuery: actical A	Technologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 CSS Tools 1 Write Less, Do More 1 spects of Web Applications 1	L15 115 116 121 121 123 125 130 131 L37
6 7	 Fro: 6.1 6.2 6.3 Pra 7.1 	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4 jQuery: actical A Web Ap	Technologies 1 Sic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 Ing Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 CSS Tools 1 Write Less, Do More 1 spects of Web Applications 1 oplications: Recap 1	L15 115 116 121 121 123 125 130 131 L37 137
6	 Fro: 6.1 6.2 6.3 Pra 7.1 7.2 	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4 jQuery: ctical A Web Ap Maintai	Technologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 A small but useful Fragment of CSS 1 CSS Tools 1 Write Less, Do More 1 oplications: Recap 1 ning State in Web Sites 1	L15 115 116 121 123 125 130 131 L37 137 140
6	 Fro: 6.1 6.2 6.3 Pra 7.1 7.2 7.3 	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4 jQuery: ctical A Web Ap Maintai Access	Technologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 A small but useful Fragment of CSS 1 Vrite Less, Do More 1 spects of Web Applications 1 oplications: Recap 1 ning State in Web Sites 1 Control and Management 1	 115 115 116 121 121 123 125 130 131 137 137 140 141
6	 From 6.1 6.2 6.3 Pra 7.1 7.2 7.3 7.4 	nt-end 2 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4 jQuery: ctical A Web A _I Maintai Access HTTPS	Technologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 A small but useful Fragment of CSS 1 CSS Tools 1 Write Less, Do More 1 oplications: Recap 1 ning State in Web Sites 1 Control and Management 1 Secure/Encrypted HTTP 1	115 115 116 121 121 123 125 130 131 137 140 141 144
6 7 8	 From 6.1 6.2 6.3 Prat 7.1 7.2 7.3 7.4 Wth 	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4 jQuery: ctical A Web A _I Maintai Access HTTPS	Technologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 A small but useful Fragment of CSS 1 CSS Tools 1 Write Less, Do More 1 oplications: Recap 1 ning State in Web Sites 1 Control and Management 1 K: Secure/Encrypted HTTP 1 ve learn in IWGS-1? 1	115 115 116 121 121 121 121 123 125 130 131 137 140 141 144 147
6 7 8 A	 From 6.1 6.2 6.3 Prat 7.1 7.2 7.3 7.4 When Exercise 	nt-end 7 Dynam 6.1.1 Cascadi 6.2.1 6.2.2 6.2.3 6.2.4 jQuery: actical A Web A _I Maintai Access HTTPS at did v	Iechnologies 1 ic HTML: Client-side Manipulation of HTML Documents 1 JavaScript in HTML 1 ng Stylesheets 1 Separating Content from Layout 1 Worked Example: The Contact Form 1 A small but useful Fragment of CSS 1 CSS Tools 1 Write Less, Do More 1 spects of Web Applications 1 oplications: Recap 1 Control and Management 1 Secure/Encrypted HTTP 1 ve learn in IWGS-1? 1	<pre>115 115 116 121 121 123 125 130 131 137 140 141 144 147 151</pre>

Informatische Werkzeuge in den Geistes- und Sozialwissenschaften 1/2 IWGSsubtitle

Prof. Dr. Michael Kohlhase

Knowledge Representation and -Processing Computer Science, FAU Erlangen-Nürnberg https://kwarc.info/kohlhase

2025-06-05

CONTENTS

Contents

0.1 Preface

0.1.1 Course Concept

Objective: The course aims at giving students an overview over the variety of digital tools and methods at the disposal of practitioners of the humanities and social sciences, explaining their intuitions on how/why they work (the way they do). The main goal of the course is to empower students for their for the emerging discipline of "digital humanities and social sciences". In contrast to a classical course in computer science which lays the mathematical and computational foundations which will become useful in the long run, we want to introduce methods and tools that can become useful in the short term and thus generate immediate success and gratification, thus alleviating the "programming shock" (the brain stops working when in contact with computer science tools or computer scientists) common in the humanities and social sciences.

Original Context: The course "Informatische Werkzeuge in den Geistes- und Sozialwissenschaften" is a first-year, two-semester course in the bachelor program "Digitale Geistes- und Sozialwissenschaften" (Digital Humanities and Social Sciences: DigiHumS) at FAU Erlangen-Nürnberg.

Open to External Students: Other bachelor programs are increasingly co-opting the course as specialization option or a key skill. There is no inherent restriction to DHSS students in this course.

Prerequisites: There are no formal prerequisites – after all it starts in the first semester for DigiHumS – but a good deal of motivation, openness towards exploring the weird and wonderful world of digital methods and tools, and a certain perseverance in the face of not understanding directly help tremendously and helps having fun in this course.

We do assume that students have a personal laptop, or access to a computer where they have admin rights, i.e. can install software. This is necessary for solving the homework. In particular, smartphones and most tablet computers will not suffice.

0.1.2 Course Contents

The course comprises two parts that are given as two-hour/week lectures.

IWGS 1 (the first semester): begins with an introduction to programming in Python which we will use as the main computational tool in the course; see ??? and ???. In particular we will cover

- systematics and culture of programming
- program and control structures
- basic data structures like numbers and strings, in particular character encodings, Unicode, and regular expressions.

Building on this, we will cover

- 1. digital documents and document processing, in particular; text files, markupsystems, HTML, and XML; see chapter 4.
- 2. basic concepts of the World Wide Web; see ???
- 3. Web technologies for interactive documents and their applications; in particular internet infrastructure, web browsers and servers, PHP, dynamic HTML, JavaScript, and CSS; see ???.

IWGS 2 (the second semester): covers selected topics and exemplary tools that will become useful in the DH. We will cover

- 1. Databases; in particular entity relationship diagrams, CRUD operations, and querying; see ???.
- 2. Image processing tools, see ???

3. Using the ontologies and the semantic web for Cultural Heritage; see ???

4. The WissKI System: A Virtual Research Environment for Cultural Heritage; see ???

5. Copyright and Data Privacy as legal foundations of DH tools; see ???

Idea: The first semester lays the foundations by introducing programming in Python and work our way towards web applications, which form the base of most modern tools in the DH. In ???, we pull all parts together to build a first, simple web application with persistent storage that manages a set of books.

After an excursion into project management systems, we introduce methods and tools for their management. Here, we extend our web application to deal with image fragments; actually building a simple replacement for a prominent DH web application.

Finally, after another excursion – this time into the legal foundations of intellectual property and data privacy the course culminates in an introduction of the WissKI system, a virtual research environment for documenting cultural heritage artifacts. Indeed the WissKI system combines all topics in the course so far.

0.1.3 Programming Exercises and JuptyterLab as a Web IDE

Programming Exercises: Most of the computer tools introduced in this course require programming e.g. for configuration, extension, or input preprocessing or work much better when the user understands the basic underlying concepts at the program level. Therefore we accompany the course with a set of (programming) exercises (given as homework to the IWGS students) that allow practicing that.

Web IDEs: In the IWGS course at FAU, which is adressed to students from the humanities and social sciences, we do not have access to a pool of standardized hardware. Students have to use their own computing devices for the programming exercises. In any group with diverse hardware, installing software, standardizing software versions, ... becomes a serious problem, even if the group only has 50 members; in IWGS, we need the Python interpreter, a text editor or integrated development environment (IDE), and various Python libraries. In IWGS we solve this by using a web IDE, which only presupposes a web browser on student hardware.

Jupyterlab: After experimenting with commercial web IDEs we settled on jupyterLab, even though it does not focus on IDE features. Jupyter notebooks allow to mix documentation, code snippets, and exercise text of programming exercises and package them into learning objects that can be downloaded, interacted with, and submitted easily. jupyterLab acts as the user interface for managing and editing jupyter notebooks and supplies standardized shell and Python REPLs for students. The jupyterLab server runs as a virtual machine on the instructor's hardware. Resource consumption is minimal in our experience (except in the week before the exam). See [JKI] for a documentation of how to set up a server for a small course like IWGS.

Limitations of JupyterLab: Of course, students who want to engage in more serious software development will eventually have to "graduate" to a regular IDE when programs become larger and more long-lived. But this – and the necessary software engineering skills – is emphatically not the focus of the IWGS course.

Exercise Notebooks: The exercise notebooks (in notebook format and PDF – unfortunately only in German) can be found at https://kwarc.info/teaching/IWGS/NB. They comprise

- outright programming exercises that introduce the Python language or allow to play with the respective concepts in Python
- code reading/debugging exercises where the character of Beatrice Beispiel almost solves interesting problems, and
- development steps towards larger applications, which often involve completing Python skeletons using the concepts taught in the lectures.

In all cases, the necessary increments to be supplied by the students are designed to not let the Python skills become a barrier, but give students the opportunity to develop the necessary programming skills in passing.

We have themed the exercises with DigiHumS topics to keep them interesting for our students.

0.1.4 This Document

Presentation: The document mixes the slides presented in class with comments of the instructor to give students a more complete background reference. **Licensing:** This document is licensed under a Creative Commons license that requires attribution, allows commercial use, and allows derivative works as long as these are licensed under the same license. **Knowledge Representation Experiment:** This document is also an experiment in knowledge representation. Under the hood, it uses the STEX package [Koh08; sTeX], a TEX/LATEX extension for semantic markup, which allows to export the contents into active documents that adapt to the reader and can be instrumented with services based on the explicitly represented meaning of the documents.

Other Resources: The lecture notes will be complemented by a selection of problems (with and without solutions) that can be used for self-study; see http://kwarc.info/teaching/IWGS.

0.1.5 Acknowledgments

Materials: The materials in this course are partially based on various lectures the author has given at Jacobs University Bremen in the years 2010-2016, these in turn have been partially based on materials and courses by Dr. Heinrich Stamerjohanns, PD Dr. Florian Rabe, and Prof. Dr. Peter Baumann. ??? have been provided by Philipp Kurth and Dr. Frank Bauer.

All course materials have been restructured and semantically annotated in the ST_EX format, so that we can base additional semantic services on them.

Teaching Assistants: The organization and material choice in the IWGS has significantly been influenced by Jonas Betzendahl and Philipp Kurth, who have been very active and dedicated teaching assistants and have given feedback on all aspects of the course. They have also provided almost all of the IWGS exercises – see subsection 0.1.3.

DigiHumS Administrators: Jacqueline Klusik-Eckert and Philipp Kurth who used to administrate the DigiHumS major at FAU together have been helpful in navigating the administrative waters of an unfamiliar faculty.

WissKI Specialists and Colleagues: ??? has profited from discussions with Peggy Große and Juliane Hamisch, then two WissKI specialists at FAU. My colleagues Prof. Peter Bell has provided the idea and data for the "Kirmes Pictures Project" that grounds some of the second semester.

JupyterLab: The JupyterLab server at https://juptyter.kwarc.info (see ???) has been developed, operated, and maintained by Jonas Betzendahl. For details see [JKI].

IWGS Students: The following students have submitted corrections and suggestions to this and earlier versions of the notes: Paul Moritz Wegener, Michael Gräwe.

0.2 Recorded Syllabus

The recorded syllabus - a record the progress of the course in the 1 - is in the course page in the ALEA system at https://courses.voll-ki.fau.de/course-home/iwgs-1. The table of contents in the IWGS lecture notes at https://kwarc.info/teaching/IWGS indicates the material covered to date in yellow.

CONTENTS

Chapter 1

Preliminaries

1.1 Administrativa

We will now go through the ground rules for the course. This is a kind of a social contract between the instructor and the students. Both have to keep their side of the deal to make learning as efficient and painless as possible.

Prerequisites		
General Prerequisites: nothing else! We will teach	Motivation, interest, cu 1 you all you need to kr	riosity, hard work. now
\triangleright You can do this course if	you want!	(we will help)
FAU	1	2025-06-05

Now we come to a topic that is always interesting to the students: the grading scheme: The short story is that things are complicated. We have to strike a good balance between what is didactically useful and what is allowed by Bavarian law and the FAU rules.





Homework assignments, quizzes, and end-semester exam may seem like a lot of work – and indeed they are – but you will need practice (getting your hands dirty) to master the IWGS concepts. We will go into the details next.

Preparedness Quizzes		
 PrepQuizzes: Before every lecture about the material from the previous starts in week 2) 	e we offer a 10 r ous week. (\sim 1	nin online quiz – the PrepQuiz 6:0?-16:15 (check on ALEA);
▷ Motivations: We do this to		
 ▷ keep you prepared and working of ▷ bonus points if the exam has ≥ ▷ update the ALEA learner mode 	continuously. 50% points I.	(primary) (potential part of your grade) (fringe benefit)
The prepquizes will be given in the bhttps://courses.ve/	ALEA system	/quiz-dash/iwgs-1
⊳ You have to be logged	d into ALEA!	(via FAU IDM)
⊳ You can take the prep	oquiz on your lap	top or phone,
$ ho \dots$ in the lecture or at	home	
$ ho \dots$ via WLAN or 4G N	letwork.	(do not overload)
▷ Prepquizzes will only be a second seco	be available ~ 16	5:0?-16:15 (check on ALEA)!
FAU	3	2025-06-05

1.2 Getting Most out of IWGS

In this section we will discuss a couple of measures that students may want to consider to get most out of the IWGS course.

None of the things discussed in this section – homeworks, tutorials, study groups, and attendance – are mandatory (we cannot force you to do them; we offer them to you as learning opportunities), but most of them are very clearly correlated with success (i.e. passing the exam and getting a good grade), so taking advantage of them may be in your own interest.



It is very well-established experience that without doing the homework assignments (or something similar) on your own, you will not master the concepts, you will not even be able to ask sensible questions, and take very little home from the course. Just sitting in the course and nodding is not enough! If you have questions please make sure you discuss them with the instructor, the teaching assistants, or your fellow students. There are three sensible venues for such discussions: online in the lectures, in the tutorials, which we discuss now, or in the course forum – see below. Finally, it is always a very good idea to form study groups with your friends.



▷ Weekly tutorials and homework assignments

(first one in week two)

⊳	Tutor: ▷ Dirk Böhme: di They know what the you learn!	(Master S rk.boehme@fau.de y are doing and really (dedie	tudent in CS) y want to help cated to DH)	
⊳ Dirk will a (grade-rele	also grade the homewor evant)	k assignments for th	e DFÜ students.	
▷ Goal 1: concept)	Reinforce what was tau	ight in class (important pillar of the IV	VGS
⊳ Goal 2:	Let you experiment wit	h Python (think of	them as Programming La	abs)
⊳ Life-savin lecture not	ng Advice: go to you tes and the homework a	r tutorial, <mark>and prepa</mark> assignments	re it by having looked at	the
⊳ Inverted	Classroom: the late	est craze in didactics	(works well if done rig	ght)
in IWGS:	lecture + homework as	ssignments + tutoria	Is $\hat{=}$ inverted classroom	
in IWGS:	lecture + homework as	ssignments + tutoria	$ls \cong inverted classroom$	
in IWGS:	lecture + homework as	ssignments + tutoria	$ls \cong inverted classroom$	
Collaborat ▷ Definition acting toge for selfish and benefit	ion n 1.2.1. Collaboration ether for common, mut benefit. In a collaborat ts from the contributio	(or cooperation) is th tual benefit, as oppo ion, every agent con ns of others.	Is $\widehat{=}$ inverted classroom 2025-06-05 The process of groups of ago sed to acting in competi tributes to the common g	ents tion goal
Collaborat Collaborat ▷ Definition acting toge for selfish and benefit ▷ In learning	ion n 1.2.1. Collaboration ether for common, mut benefit. In a collaborat ts from the contributio g situations, the benefit	signments + tutoria (or cooperation) is the tual benefit, as oppo- tion, every agent con ns of others. is "better learning".	Is $\widehat{=}$ inverted classroom 2025-06-05 The process of groups of ago used to acting in competi- tributes to the common g	ents tion goal
In IWGS: Collaborat ▷ Definition acting toge for selfish and benefit ▷ In learning ▷ Observat than in con	ion ion n 1.2.1. Collaboration ether for common, mut benefit. In a collaborat ts from the contributio g situations, the benefit ion: In collaborative le mpetitive learning.	ssignments + tutoria (or cooperation) is the tual benefit, as oppo- cion, every agent con ns of others. : is "better learning". arning, the overall res	Is $\widehat{=}$ inverted classroom 2025-06-05 The process of groups of agoresed to acting in competi- tributes to the common groups of th	ents tion goal

1. A Those learners who work/help most, learn most!
 2. Freeloaders – individuals who only watch – learn very little!

\triangleright	It is OK to collaborate on homewor	k a	ssignments in IWGS!	(no bonu	ıs points)
\triangleright	Choose your study group well!		(ALeA helps via the s	study buddy	/ feature)
		6		2025.06.05	

As we said above, almost all of the components of the IWGS course are optional. That even applies to attendance. But make no mistake, attendance is important to most of you. Let me explain, ...



1.3. LEARNING RESOURCES FOR



1.3 Learning Resources for IWGS



FAU has issued a very insightful guide on using lecture videos. It is a good idea to heed these recommendations, even if they seem annoying at first.





NOT a Resource for : LLMs – AI-tools like ChatGPT

- ▷ Definition 1.3.1. A large language model (LLM) is a computational model capable of language generation or other natural language processing tasks.
- ▷ **Example 1.3.2.** OpenAI's GPT, Google's Bard, and Meta's Llama.
- Definition 1.3.3. A chatbot is a software application or web interface that is designed to mimic human conversation through text or voice interactions. Modern chatbots are usually based on LLMs.
- ▷ Example 1.3.4 (ChatGPT talks about IWGS). (Aha, where does this come from?)





1.4 Goals, Culture, & Outline of the Course



One of the most important tasks in an inter/trans-disciplinary enterprise – and that what "digital humanities" is, fundamentally – is to understand the disciplinary language, intuitions and foundational assumptions of the respective other side. Assuming that most students are more versed in the "humanities and social sciences" side we want to try to give an overview of the "computer science culture".

Academic Culture in Computer Science

> Definition 1.4.1. The academic culture is the overall style of working, research,

and discussion in an academic field.				
Observation 1.4.2. There are significant differences in the academic culture be- tween computer science, the humanities and the social sciences.				
▷ Computer science is an engineering discipline (we build the				
▷ given a problem we look for a (mathematic	al) model, we can think with			
\triangleright once we have one, we try to re-express it w	ith fewer "primitives" (concepts)			
▷ once we have, we generalize it (make it more widely applicabl				
\triangleright only then do we implement it in a program (ideally				
Design of versatile, usable, and elegant tools is an important concern				
▷ Almost all technical literature is in English. (technical vocabulary too)				
▷ CSlings love shallow hierarchies. (no personality cult; alle per Du				
FAU : 12 2025-06-05 EVENERAL				

Please keep in mind that – self-awareness is always difficult – the list above may be incomplete and clouded by mirror-gazing. We now come to the concrete topics we want to cover in IWGS. The guiding intuition for the selection is to concentrate on techniques that may become useful in day-to-day DH work – not CS completeness or teaching efficiency.



1.5 ALeA – AI-Supported Learning

In this section we introduce the ALEA (Adaptive Learning Assistant) system, a learning support

1.5. ALEA – AI-SUPPORTED LEARNING

system we will use to support students in IWGS.

ALEA: Adaptive Learn	ing Assistant	
▷ Idea: Use AI methods to h	nelp teach/learn Al	(AI4AI)
Concretely: Provide HTM learning support services into	L versions of the IWGS sl o them. (for	ides/lecture notes and embed pre/postparation of lectures)
Definition 1.5.1. Call a do information needs of the rea	cument <mark>active</mark> , iff it is inte ders.	eractive and adapts to specific (lecture notes on steroids)
\triangleright Intuition: ALEA serves a	ctive course materials.	(PDF mostly inactive)
\triangleright Goal: Make ALEA more l	ike a instructor + study g	group than like a book!
► Example 1.5.2 (Course N	otes).	Image: State Programming Features) Significant: This will be able to grant and the provide to the search tree reture. Image: State Programming Features) State Programming Features) Image: State Programming Features Image: State Programing Features <td< td=""></td<>
ightarrow yellow parts in table of co	ontents (left) already cove	ered in lectures.
FAU	14	2025-06-05

The central idea in the AI4AI approach – using AI to support learning AI – and thus the ALeA system is that we want to make course materials – i.e. what we give to students for preparing and postparing lectures – more like teachers and study groups (only available 24/7) than like static books.





The ALEA IWGS page is the central entry point for working with the ALEA system. You can get to all the components of the system, including two presentations of the course contents (notesand slides-centric ones), the flashcards, the localized forum, and the quiz dashboard.

We now come to the heart of the ALeA system: its learning support services, which we will now briefly introduce. Note that this presentation is not really sufficient to undertstand what you may be getting out of them, you will have to try them, and interact with them sufficiently that the learner model can get a good estimate of your competencies to adapt the results to you.



Example 1.5.4 (More Definitions on Click). Clicking on a (cyan) term reference shows us more definitions from other contexts.

1.5. ALEA – AI-SUPPORTED LEARNING





Note that this is only an initial collection of learning support services, we are constantly working on additional ones. Look out for feature notifications ($\bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc$) on the upper right hand of the ALeA screen.

(Practice/Remedial) Problems Everywhere				
▶ Problem: Learning requires a mix of understanding and test-driven practice.				
Idea: ALeA supplies targeted practice problems everywhere.				
Concretely: Revision markers at the end of sections.				
▷ A relatively non-intrusive overview over competency				
Review Minimax Search V				
\triangleright Click to extend it for details.				
Review Minimax Search				
PRACTICE PROBLEMS (7)				
▷ Practice problems as usual. (targeted to your specific competency)				

	Review Minimax Search	^	
	Problem 6 of 7		
	(Minimax) which of the following statements about minimax are true?		
	An extension a or are using function a to inner nodes, a is computed recursively. Max attempts to maximize $\hat{u}(s)$ of states reachable during play.		
	Minimax computes an online strategy Returns an optimal action, assuming perfect opponent play		
	CHECK SOLUTION		
FAU :	17 202	5-06-05	

While the learning support services up to now have been adressed to individual learners, we now turn to services addressed to communities of learners, ranging from study groups with three learners, to whole courses, and even – eventually – all the alumni of a course, if they have not de-registered from ALeA.

Currently, the community aspect of ALeA only consists in localized interactions with the course materials.

The ALeA system uses the semantic structure of the course materials to localize some interactions that are otherwise often from separate applications. Here we see two:

- 1. one for reporting content errors and thus making the material better for all learners and ''
- 2. a localized course forum, where forum threads can be attached to learning objects.

Localized Interactions with the Community				
 Selecting text brings up localized – i.e. anchored on the selection – interactions: post a (public) comment or take (private) note report an error to the course authors/instructors 				
Localized comments induce a thread in the ALEA forum (like the StudOn Forum, but targeted towards specific learning objects.) problem in the abstract, i.e. make a plan before we actually enter the situation (i.e. offine), and then when the problem arise, only one grant thread out thread out the situation (i.e. offine), and then when the problem arise, only one grant thread out thread out thread out thread out the situation (i.e. offine).				
situation a MY NOTES COMMENTS r the actions of others). As this is much more				
Pro 1 comments C Y				
A sequence of actions is a solution ▷ In o It could equivalently be defined as a sequence of actions: we can compute the state sequence from the action sequence and – given the initial state – the action sequence for the state sequence. a chance to find general algorithms.				
▷ Con Request response ▷ S ▷ A POST				
A sec Michael Kohlhase ① 4 minutes ago ▲ REPLY : pail state. Problem solving computes solutions from (A sequence of actions is a solution)				
Defi Too not understan this, why is it a solution a sequence or states? sequence based complete knowledge of the envir CLOSE				
 Definition 1.1.3. In online problem solving an agent computes one action at a time based on incoming perceptions. 				



We can use the same four models discussed in the space of guided tours to deploy additional learning support services, which we now discuss.

New Feature: Drilling with Flas	hcards
ightarrow Flashcards challenge you with a task (1	term/problem) on the front
🗱 COURSES 🔺 🔹 🗮 Michael	COURSES 🔺 🔹 🕫 🚟 Michael
	Definition 0.1. The weight space is the space of all possible combinations of weights. Loss minimization in a weight space is called weight fitting.
weight space	
Assess Your Comptence:	Assess Your Comptence:
and the definition/answer is on the	back.
\triangleright Self-assessment updates the learner mo	odel (before/after)
Idea: Challenge yourself to a card st the learner model eliminates all.	ack, keep drilling/assessing flashcards until
Bonus: Flashcards can be generated equivalent to free beer)	from existing semantic markup (educational
FAU : 19	2025-06-05

We have already seen above how the learner model can drive the drilling with flashcards. It can also be used for the configuration of card stacks by configuring a domain e.g. a section in the course materials and a competency threshold. We now come to a very important issue that we always face when we do AI systems that interface with humans. Most web technology companies that take one the approach "the user pays for the services with their personal data, which is sold on" or integrate advertising for renumeration. Both are not acceptable in university setting.

But abstaining from monetizing personal data still leaves the problem how to protect it from intentional or accidental misuse. Even though the GDPR has quite extensive exceptions for research, the ALeA system – a research prototype – adheres to the principles and mandates of the GDPR. In particular it makes sure that personal data of the learners is only used in learning support services directly or indirectly initiated by the learners themselves.

Learner Data and Privacy in ALEA

1.5. ALEA – AI-SUPPORTED LEARNING

▷ Observation: Learning support	services in ALEA use the learner model; they
\triangleright need the learner model data t	o adapt to the invidivual learner!
▷ collect learner interaction data	a (to update the learner model)
Consequence: You need to be learning support services!	ogged in (via your FAU IDM credentials) for useful
▷ Problem: Learner model data is	s highly sensitive personal data!
▷ ALeA Promise: The ALEA to safe.	eam does the utmost to keep your personal data (SSO via FAU IDM/eduGAIN, ALEA trust zone)
▷ ALeA Privacy Axioms:	
1. $ALEA$ only collects learner mo	odels data about logged in users.
2. Personally identifiable learner n (delegation possible)	nodel data is only accessible to its subject
3. Learners can always query the	learner model about its data.
 All learner model data can be usability deterioration) 	e purged without negative consequences (except
5. Logging into $ALEA$ is complet	ely optional.
Observation: Authentication for you can always purge the learner	or bonus quizzes are somewhat less optional, but model later.
FAU	20 2025-06-05 CONTRACTOR

So, now that you have an overview over what the ALEA system can do for you, let us see what you have to concretely do to be able to use it.



16				CHAPTER 1.	PRELIMINARIES
L	Fau	:	21	2025-06-05	

Even if you did not understand some of the AI jargon or the underlying methods (yet), you should be good to go for using the ALEA system in your day-to-day work.

Chapter 2

Introduction to Programming

What is Programming? 2.1

Programming is an important and distinctive part of "Informatische Werkzeuge in den Geistesund Sozialwissenschaften" – the topic of this course. Before we delve into learning Python, we will review some of the basics of computing to situate the discussion.

To understand programming, it is important to realize that computers are universal machines. Unlike a conventional tool e.g a spade – which has a limited number of purposes/behaviors – digging holes in case of a spade, maybe hitting someone over the head, a computer can be given arbitrary¹ purposes/behaviors by specifying them in form of a program.

This notion of a program as a behavior specification for an universal machine is so powerful, that the field of computer science is centered around studying it – and what we can do with programs, this includes

- i) storing and manipulating data about the world,
- ii) encoding, generating, and interpreting image, audio, and video,
- *iii*) transporting information for communication,
- iv) representing knowledge and reasoning,
- v) transforming, optimizing, and verifying other programs,
- vi) learning patterns in data and predicting the future from the past.

Computer Hardware/Software & Programming

Definition 2.1.1. Computers consist of hardware and software.
 Definition 2.1.2. Hardware consists of

¹as long as they are "computable", not all are.



A universal machine has to have – so experience in computer science shows certain distinctive parts.

- A CPU that consists of a
 - control unit that interprets the program and controls the flow of instructions and
 - a arithmetic/logic unit (ALU) that does the actual computations internally.
- Memory that allows the system to store data during runtime (volatile storage; usually RAM) and between runs of the system (persistant storage; usually hard disks, solid state disks, magnetic tapes, or optical media).
- I/O devices for the communication with the user and other computers.

With these components we can build various kinds of universal machines; these range from thought experiments like Turing machines, to today's general purpose computers like your laptop with various embedded systems (wristwatches, Internet routers, airbag controllers, ...) in-between. Note that – given enough fantasy – the human brain has the same components. Indeed the human mind is a universal machine – we can think whatever we want, react to the environment, and are not limited to particular behaviors. There is a sub-field of computer science that studies this: AI (artificial intelligence). In this analogy, the brain is the "hardware" –sometimes called "wetware" because it is not made of hard silicon or "meat machine"². It is instructional to think about what the program and the data might be in this analogy.



 \triangleright completely stupid: will not do what you mean unless you tell it exactly

²Marvin Minsky; one of the founding fathers of the field of artificial intelligence

2.1. WHAT IS PROGRAMMING?



AI studies human intelligence with the premise that the brain is a computational machine and that intelligence is a "program" running on it. In particular, the working hypothesis is that we can "program" intelligence. Even though AI has many successful applications, it has not succeeded in creating a machine that exhibits the equivalent to general human intelligence, so the jury is still out whether the AI hypothesis is true or not. In any case it is a fascinating area of scientific inquiry.

Note: This has an immediate consequence for the discussion in our course. Even though computers can execute programs very efficiently, you should not expect them to "think" like a human. In particular, they will execute programs exactly as you have written them. This has two consequences:

- the behavior of programs is in principle predictable
- all errors of program behavior are your own (the programmer's)

In computer science, we distinguish two levels on which we can talk about programs. The more general is the level of algorithms, which is independent of the concrete programming language. Algorithms express the general ideas and flow of computation and can be realized in various languages, but are all equivalent – in terms of the algorithms they implement.

As they are not bound to programming languages algorithms transcend them, and we can find them in our daily lives, e.g. as sequences of instructions like recipes, game instructions, and the like. This should make algorithms quite familiar; the only difference of programs is that they are written down in an unambiguous syntax that a computer can understand.





We have two kinds of programming languages: one which the CPU can execute directly – these are very very difficult for humans to understand and maintain – and higher-level ones that are understandable by humans. If we want to use high-level languages – and we do, then we need to have some way bridging the language gap: this is what compilers and interpreters do.

2.2 Programming in IWGS

After the general introduction to programming in ???, we now instantiate the situation to the IWGS course, where we use Python as the primary programming language.



Note that IWGS is not a programming course, which concentrates on teaching a programming language in all it gory detail. Instead we want to use the IWGS lectures to introduce the necessary concepts and use the tutorials to introduce additional language features based on these.



However, the result would probably be the following:



If we just start hacking before we fully understand the problem, chances are very good that we will waste time going down blind alleys, and garden paths, instead of attacking problems. So the main motto of this course is:



2.3 Programming in Python

In this section we will introduce the basics of the Python language. Python will be used as our means to express algorithms and to explore the computational properties of the objects we introduce in IWGS.

2.3.1 Hello IWGS

Before we get into the syntax and meaning of Python, let us recap why we chose this particular language for IWGS.



Installing Python: Python can be installed from http://python.org \sim "Downloads", as a MSWindows installer or a macOS disk image. For linux it is best installed via the package manager, e.g. using

2.3. PROGRAMMING IN PYTHON

```
sudo apt—get update
```

sudo apt—get install python3.7

The download will install the Python interpreter and the Python shell IDLE3 that can be used for interacting with the interpreter directly.

It is important that you make sure (tick the box in the Windows installer) that the python executable is added to the path. In the shell¹, you can then use the command

```
python (filename)
```

to run the python file $\langle\!\langle$ filename $\rangle\!\rangle$. This is better than using the windows-specific

```
py ((filename))
```

which does not need the python interpreter on the path as we will see later.

Arithmetic Expressions in Pytl	hon	
▷ Expressions are "programs" that com	pute values	(here: numbers)
⊳ Integers (numb	ers without a decimal point)	
 ▷ operators: addition (+ tion (*), division (/), i der/modulo (%), ▷ Division yields a float 	-), subtraction (), multiplica- integer division (//), remain-	Python 3.1.3 [GCC 4.5.1 201] Type "copyrigh"
⊳ Floats (nu	mbers with a decimal point)	7 >>> 3 - 4
 Operators: integer b (ceil), exponential (exp 	elow (floor), integer above p), square root (sqrt),	-1 >>> 3 - 4.0 -1.0
 Numbers are values, i.e computed with. (reference _) 	. data objects that can be the last computed one with	>>> 3 * 4 12 >>> 27 / 5 5.4 >>> 27 // 5
Definition 2.3.1. Express (and other expressions) vi	sions are created from values a operators.	5 >>> 27 % 5 2
Observation: The Pyth pressions to values by con	non interpreter simplifies ex- nputation.	>>>
FAU .	30 202	5-06-05 COMBINEST

Before we go on to learn more basic Python operators and instructions, we address an important general topic: comments in program code.

 Comments in Python

 Generally: It is highly advisable to insert comments into your programs,
 especially, if others are going to read your code,
 you may very well be one of the "others" yourself,
 (in a year's time)
 writing comments first helps you organize your thoughts.

 $^{^1\}mathrm{EdNOTE:}$ fully introduce the concept of a shell in the next round



2.3.2 JupyterLab, a Python Web IDE for IWGS

In IWGS, we want to use the jupyterLab cloud service. This runs the Python interpreter on a cloud server and gives you a browser window with a web IDE, which you can use for interacting with the interpreter. You will have to make an account there; details to follow.

jupyterLab A Cloud IDI	E for Python		
Eor belging you it would be	$rac{1}{1}$	could access to your code	
		could access to your code	
Idea: Use a web IDE (a we Lab, which you can use for i	b based integrated interacting with th	l development environment): le interpreter.	jupyter-
\triangleright We will use jupyterLab for I	IWGS.	but you can also use Pytho	n locally)
▷ Homework: Set up jupyte	erLab		
\triangleright make an account at htt	p://jupyter.kw	arc.info	
FAU	32	2025-06-05	COMB RIGHTS RESERVED

The advantage of a cloud IDE like jupyterLab for a course like IWGS is that you do not need any installation, cannot lose your files, and your teachers (the course instructor and the teaching assistants) can see (and even directly interact with) the your run time environment. This gives us a much more controlled setting and we can help you better.

Both IDLE3 as well as jupyterLab come with an integrated editor for writing Python programs. These editors gives you Python syntax highlighting, and interpreter and debugger integration. In short, IDLE3 and jupyterLab are integrated development environments for Python. Let us now go through the interface of the jupyterLab IDE.

jupyterLab Components

▷ **Definition 2.3.2.** The jupyterLab dashboard gives you access to all components.

2.3. PROGRAMMING IN PYTHON



	C JupyterLab	× +		
$(\leftarrow) \rightarrow $ G \cdot	۵ () ۳۰	https://jupyter.kwarc.in ••• 🗟	אל ≫	=
Getting Started	d 📄 FAU 📄 Ser	vices 📄 News 📄 MathWeb 📄 Jacobs	AG Lists	»
	view Run	Rernel labs Settings Help		
	0	mkoninase@jupyter: ~ ×		_
Name	•	<pre>mkohlhase@jupyter:~\$ pwd /home/mkohlhase mkohlhase@jupyter:~\$ ls</pre>		
🗖 Untitle	d.ipynb	test.py Untitled.ipynb mkohlhase@jupyter:~\$		
🕐 📿 test.py	/	_		
Definition 2.3.5. A a computer's operation There are multiple shows	shell is a c ng system. ell impleme	ommand line interface for ntations: sh, csh, bash, zs	accessing	the services of er in advanced
features.				
⊳ Useful shell comma	ands: See	e.g. [All18] for a basic tu	torial	
\triangleright ls: "list" the files	in this dire	ctory		
⊳ mkdir: "make" fo	lder (called	"directory")		
⊳ pwd: "print worki	ng director	y"		(where am I)
⊳ cd 《dirname》: "c	hange dire	ctory"		,
⊳ if ∥dirname»	= : one u	n in the directory tree		
⊳ empty ∛dirnai	me»: go to	your home directory.		
⊳ rm ∥name\\: rem	ove file/dir	ectory		
⊳ cn/my //filenama		noll: convite or rename		
⊳ cp/mv ((mename)	// ((newnam	. copy to or rename		
⊳ cp/mv ((filename)	» «dirname	e»: copy or move to		
⊳ see [All18] for	more			
FAU		33	2025-06	-05 SYMMERSENSES

Now that we understand our tools, we can wrote our first program: Traditionally, this is a "hello-world program" (see [HWC] for a description and a list of hello world programs in hundreds of languages) which just prints the string "Hello World" to the console. For Python, this is very simple as we can see below. We use this program to explain the concept of a program as a (text) file, which can be started from the console.

A first program in Python > A classic "Hello World" program: start your python console, type print("Hello IWGS").
(print a string)

.

	Consolution Conso	× + ∞ /jupyter.kwarc.info/user/ ···· ♡ ☆ III © Ø Ø Ξ eves MathWeb jacobs A 6 bis Settings Help
	Name Vytho Type Dutitide Folder Uutitide folder Vytho Vy	on 3.5.3 (default, Sep 27 2018, 17:25:39) 'copyright', 'credits' or 'license' for more information non 7.8.0 — An enhanced Interactive Python. Type '7' for help. I: print("Hello IWGS") Hello IWGS
	1 💽 1 🐵 Python 3 Idle	11 Ln 1, Col 1 Console 2
1. got to 2. Type y	the jupyterLab dashboa our program, Ċ File Edit View Run Ke	ard select "Text File", ernel Tabs Settings Help
	+ b t C / Name Erste_Schritte_mit_Ju hello.py	<pre> E hello.py × 1 # my first python program print("Hello IWGS") </pre>
(test.py Untitled.ipvnb 	
ſ	 test.py Untitled.ipynb Untitled.ipynb S 3 Python 	Ln 2, Col 21 Spaces: 4 hello.py
3. Save th 4. Go to y	 test.py Untitled.ipynb 3 [®] Python The file as hello.py Your terminal and type 	Ln 2, Col 21 Spaces: 4 hello.py
3. Save th 4. Go to y 3' Altern	 test.py Untitled.ipynb 3 Python re file as hello.py rour terminal and type atively: go to your python 	Ln 2, Col 21 Spaces: 4 hello.py python3 hello.py thon console and type (in the same directory)
3. Save th 4. Go to y 3' Altern import	 test,py Untitled.ipynb 3 Python as hello.py your terminal and type atively: go to your pythe 	Ln 2, Col 21 Spaces: 4 hello.py python3 hello.py thon console and type (in the same directory)

We have seen that we can just call a program from the terminal, if we stored it in a file. In fact, we can do better: we can make our program behave like a native shell instruction.

- 1. The file extension .py is only used by convention, we can leave it out and simply call the file hello.
- 2. Then we can add a special Python comment in the first line

python (filename)

which the terminal interprets as "call the program python3 on me".

3. Finally, we make the file hello executable, i.e. tell the terminal the file should behave like a shell command by issuing

chmod u+x booksapp

in the directory where the file hello is stored.
4. We add the line

export PATH="./:\${PATH}"

to the file .bashrc. This tells the terminal where to look for programs (here the respective current directory called .)

With this simple recipe we could in principle extend the repertoire of instructions of the terminal and automate repetitive tasks.

We now come to the signature component of jupyterLab: jupyter notebooks. They take the important practice of documenting code to a whole new level. Instead of just allowing comments in program files, they provide rich text cells, in which we can write elaborate text.

Jupyter Notebooks
▷ Definition 2.3.6. Jupyter notebooks are documents that combine live runnable code with rich, narrative text (for comments and explanations).
▷ Definition 2.3.7. Jupyter notebooks consist of cells which come in three forms:
 ▷ a raw cell shows text as is, ▷ a markdown cell interprets the contents as markdown text, (later more) ▷ a code cell interprets the contents as (e.g. Python) code.
\triangleright Cells can be executed by pressing "shift enter". (Just "enter" gives a new line)
▷ Idea: Jupyter notebooks act as a REPL, just as IDLE3, but allows
 documentation in raw and markdown cells and changing and re-executing existing cells.
FAU : 35 2025-06-05 CONTRACT

Jupyter Notebooks

▷ Example 2.3.8 (Showing off Cells in a Notebook).

28

2.3. PROGRAMMING IN PYTHON

.

Name Name Worklad Folder Untitled Folder Untitled Folder Untitled Folder Untitled Folder Untitled Joynb M Untitled Joynb M Untitled Joynb	A + K A Markdown Cell This is for explanation, we will also show computation cells raw cells and now some python [11]: print("hello IM65") hello IM65 [10]: 344 [10]: 7 ## ray cells showy ray text, e.o. for	Python 3 O	
1 🛐 3 🐵 Python 3 Idle	(): Mode: Command 🛞 Ln 1, Cr	II 4 Untitled1.ipynb Show All X	
			8

- \triangleright Idea: We can translate between markup formats.
- ▷ Definition 2.3.9. Markdown is a family of markup formats whose control words are unobtrusive and easy to write in a text editor. It is intended to be converted to HTML and other formats for display.
- ▷ **Example 2.3.10.** Markdown is used in applications that want to make user input easy and efficient, e.g. wikis and issue tracking systems.
- ▷ Workflow: Users write markdown, which is formatted to HTML and then served for display.
- A good cheet-sheet for markdown control words can be found at https://github. com/adam-p/markdown-here/wiki/Markdown-Cheatsheet.

37

2.3.3 Variables and Types

Fau

And we start with a general feature of programming languages: we can give names to values and use them multiple times. Conceptually, we are introducing shortcuts, and in reality, we are giving ourselves a way of storing values in memory so that we can reference them later.

Variables in Python
▷ Idea: Values (of expressions) can be given a name for later reference.
> Definition 2.3.11. A variable is an identifier (the variable name) that references

29

2025-06-05



Let us fortify our intuition about variables with some examples. The first shows that we sometimes need variables to store objects out of the way and the second one that we can use variables to assemble intermeditate results.

```
Variables in Python: Extended Example
 ▷ Example 2.3.14 (Swapping Variables). To exchange the values of two variables,
  we have to cache the first in an auxiliary variable.
  a = 45
  b = 0
  print("a =", a, "b =", b)
  print("Swap the contents of a and b")
  swap = a
  a = b
   b = swap
  print("a =", a, "b =", b)
  Here we see the first example of a Python script, i.e. a series of Python commands,
  that jointly perform an action (and communicates it to the user).
 Example 2.3.15 (Variables for Storing Intermediate Variables).
   >> x = "OhGott"
   >> y = x + x + x
   >> \dot{z} = y + y + y
```



If we use variables to assemble intermediate results, we can use telling names to document what these intermediate objects are – something we did not do well in Example 2.3.15; but admittely, the meaning of the objects in this contrived example is questionable.

The next phenomenon in Python is also common to many (but not all) programming languages: expressions are classified by the kind of objects their values are. Objects can be simple (i.e. of a basic type; Python has five of these) or complex, i.e. composed of other objects; we will go into that below.

Recall: Python programs process data (values), which can be combined by oper- ators and variable into expressions.					
⊳ Dat	ta types grou	up data an	d tell the interpreter w	what to expect	
	1, 2, 3, etc. "hello" is da	are data c ata of type	of type "integer" "string"		
⊳ Dat	ta types dete	ermine whi	ch operators can be a	pplied	
 In Python, every values has a type, variables can have any type, but can only be assigned values of their type. Definition 2.3.16. Python has the following five basic types 					
⊳ In F assig ⊳ Def	Python, ever gned values finition 2.3.	ry values h of their typ . 16. Pytho	has a type, variables ca pe. on has the following fiv	an have any type, but can only be ve basic types	
⊳ In F assig ⊳ Def	Python, ever gned values finition 2.3.	ry values h of their typ . 16. Pythc Keyword	has a type, variables ca pe. on has the following fiv contains	an have any type, but can only be ve basic types Examples	
⊳ In F assig ⊳ Def	Python, ever gned values finition 2.3. Data type integers	ry values h of their typ .16. Pytho Keyword int	has a type, variables cape.	 an have any type, but can only be ve basic types Examples 1, -5, 0, 	
⊳ In F assig ⊳ Def	Python, ever gned values finition 2.3. Data type integers floats	ry values h of their typ . 16. Pytho Keyword int float	as a type, variables cape. on has the following fiv contains bounded integers floating point numbers	Examples 1, -5, 0, 1.2, .125, -1.0,	
⊳ In F assi _ℓ ⊳ De	Python, ever gned values finition 2.3. Data type integers floats strings	ry values h of their typ .16. Pytho Keyword int float str	as a type, variables cape. on has the following fiv contains bounded integers floating point numbers strings	an have any type, but can only be ve basic types Examples 1, -5, 0, 1.2, .125, -1.0, "Hello", 'Hello', "123", 'a',	
⊳ In F assi _€ ⊳ Def	Python, ever gned values finition 2.3. Data type integers floats strings Booleans	ry values h of their typ .16. Pytho Keyword int float str bool	as a type, variables cape. on has the following fiv contains bounded integers floating point numbers strings truth values	an have any type, but can only be ve basic types Examples 1, -5, 0, 1.2, .125, -1.0, "Hello", 'Hello', "123", 'a', True, False	
⊳ In F assig ⊳ Def	Python, ever gned values finition 2.3. Data type integers floats strings Booleans complexes	ry values h of their typ .16. Pytho Keyword int float str bool complex	as a type, variables cape. on has the following fiv contains bounded integers floating point numbers strings truth values complex numbers	an have any type, but can only be re basic types Examples 1, -5, 0, 1.2, .125, -1.0, "Hello", 'Hello', "123", 'a', True, False 2+3j,	

We will now see what we can – and cannot – do with data types, this becomes most noticable in variable assignments which establishes a type for the variable (this cannot be change any more) and in the application of operators to arguments (which have to be of the correct type).

Data Types in Python (continued) ▷ The type of a variable is automatically determined in the first variable assignment (before that the variable is unbound) >>> firstVariable = 23 # integer >>> type(firstVariable) <class 'int'> weight = 3.45 # float first = 'Hello' # str

CHAPTER 2. INTRODUCTION TO PROGRAMMING

⊳ Hint: class bi	The Python t)	function type to computes the type	(don't worry a	about the
Fau	÷	41	2025-06-05	COME ADDITIS ADDIANAED



2.3.4 Python Control Structures

So far, we only know how to make programs that are a simple sequence of instructions no repetitions, no alternative pathways. Example 2.3.13 is a perfect example. We will now change that by introducing control structures, i.e complex program instructions that change the control flow of the program.



32

2.3. PROGRAMMING IN PYTHON



After this general introduction – conditional execution and loops are supported by all programming languages in some form – we will see how this is realized in Python



Python uses indenting to signify nesting of body parts in control structures – and other structures as we will see later. This is a very un-typical syntactic choice in programming languages, which typically use brackets, braces, or other paired delimiters to indicate nesting and give the freedom of choice in indenting to programmers. This freedom is so ingrained in programming practice, that we emphasize the difference here. The following example shows conditional execution in action.

Conditional Execution Example
\triangleright Example 2.3.25 (Empathy in Python).
answer = input("Are you happy? ")
if answer == 'No' or answer == 'no':
<pre>print("Have a chocolate!")</pre>
else:
print("Good!")
print("Can I help you with something else?")
Note the indenting of the body parts.
BTW: input is an operator that prints its argument string, waits for user input, and returns that.
FAU : 45 2025-06-05 EXTENSION

But conditional execution in Python has one more trick up its sleeve: what we can do with two branches, we can do with more as well.



Note that the **elif** is just "syntactic sugar" that does not add anything new to the language: we could have expressed the same functionality as two nested if/else statements

But this would have introduced an additional layer of nesting (per **elif** clause in the original). The nested syntax also obscures the fact that all branches are essentially equal. Now let us see the syntax for loops in Python.



Fau	:	47	2025-06-05	COMBIDENTIES ANAL

As always we will fortify our intuition with a couple of small examples.

Examples of Loops			
⊳ Example 2.3.28 (Counti	ng in python).		
# Prints out 0,1,2,3,4			
count = 0			
while count < 5 :			
print(count)			
count $+= 1 \#$ This is	the same as $count = council $	unt + 1	
This is the standard patter incrementing it in every pa Example 2.3.29 (Breaki	n for using while : using a iss through the loop. ng an unbounded Loop	loop variable (here cour).	it) and
# Prints out 0,1,2,3,4 but	uses break		
count = 0			
while True:			
print (count)			
$\operatorname{count} += 1$			
If count ≥ 5 :			
break			
FAU	48	2025-06-05	

Example 2.3.28 and Example 2.3.29 do the same thing: counting from zero to four, but using different mechanisms. This is normal in programming there is not "one correct solution". But the first solution is the "standard one", and is preferred, sind it is shorter and more readable. The **break** functionality shown off in the second one is still very useful. Take for instance the problem of computing the product of the numbers -10 to 1.000.000. The naive implementation of this is on the left below which does a lot of unnecessary work, because as soon was we passed 0, then the whole product must be zero. A more efficient implementation is on the right which breaks after seeing the first zero.

Direct Implementation

More Efficient

 $\begin{array}{l} \mbox{count} = -10 \\ \mbox{prod} = 1 \\ \mbox{while count} < 1000000: \\ \mbox{prod} \ *= \mbox{count} \\ \mbox{count} \ += 1 \end{array}$

```
count = -10
prod = 1
while count <= 1000000:
    prod *= count
    if count = 0 :
        break
        count += 1</pre>
```

Examples of Loops

▷ Example 2.3.30 (Exceptions in the Loop).

```
# Prints out only odd numbers - 1,3,5,7,9
count = 0
while count < 10
```

count	t += 1			
# Ch	neck if x is even			
if cou	unt $\% 2 == 0$:			
C	continue			
print	(count)			
Fau	:	49	2025-06-05	SUME AISHIN RESERVED

2.4 Some Thoughts about Computers and Programs

Finally, we want to go over a couple of general issues pertaining to programs and (universal) machines. We will just go over them to get the intuitions – which are central for understanding computer science and let the lecture "Theoretical Computer Science" fill in the details and justifications.

Computers as Universal Machines (a taste of theoretical CS)				
Observation: Computers are universal tools: their behavior is determined by a program; they can do anything, the program specifies.				
Context: Tools in most other disciplines are specific to particular tasks. (except in e.g. ribosomes in cell biology)				
Remark 2.4.1 (Deep Fundamental Result). There are things no computer can compute.				
Example 2.4.2. There cannot be a program that decides whether another program will terminate in finite time.				
▷ Remark 2.4.3 (Church-Turing Hypothesis). There are two classes of languages				
Turing complete (or computationally universal) ones that can compute what is theoretically possible.				
▷ data languages that cannot. (but describe data sets)				
Observation 2.4.4 (Turing Equivalence). All programming languages are (made to be) universal, so they can compute exactly the same. (compilers/interpreters exist)				
▷in particular: Everybody who tells you that one programming languages is the best has no idea what they're talking about (though differences in efficiency, convenience, and beauty exist)				
EAU : 50 2025-06-05 EXERCISE				

Artificial Intelligence

- ▷ Another Universal Tool: The human mind. (We can understand/learn anything.)
- **Strong Artificial Intelligence:** claims that the brain is just another computer.

2.4. SOME THOUGHTS ABOUT COMPUTERS AND PROGRAMS

\triangleright If that is true then			
▷ the human mind underlies▷ we may be able to find the	the same restrictions the same restriction e "mind-program".	ns as computational mach	ines
FAU	51	2025-06-05	CC)

We now come to one of the most important, but maybe least acknowledged principles of programming languages: The principle of compositionality. To fully understand it, we need to fix some fundamental vocabulary.



All of this is very abstract – it has to be as we have not fixed a programming language yet and you will only understand the true impact of the compositionality principle over time and with programming experience. Let us now see what this means concretely for our course.





2.5 More about Python

After we have had some general thoughts about programming in general, we can get back to concrete Python facilities and idoms. We will concentrate on those – there are lots and lots more – that are useful in IWGS.

2.5.1 Sequences and Iteration

We now come to a commonly used class of objects in Python: sequences, such as lists, sets, tuples, and ranges as well as dictionaries.

They are used for storing, accumulating, and accessing objects in various ways in programs. They all have in common, that they can be used for iteration, thus creating a uniform interface to similar functionality.

Lists in Python

- ▷ **Definition 2.5.1.** A list is a finite sequence of objects, its elements.
- In programming languages, lists are used for locally storing and passing around collections of objects.
- ▷ In Python lists can be written as a sequence of comma separated expressions between square brackets.
- \triangleright **Definition 2.5.2.** We call [((seq))] the list constructor.
- ▷ Example 2.5.3 (Three lists). Elements can be of different types in Python

list1 = ['physics', 'chemistry', 1997, 2000]; list2 = [1, 2, 3, 4, 5]; list3 = ["a", "b", "c", "d"];

Example 2.5.4. List elements can be accessed by specifying ranges



38

Definition 2.5.5 slicing.	. Selecting sublists by specifying	start and/or end is called	ist
FAU	55	2025-06-05)

As Example 2.5.4 shows, Python treats counting in list accessors somewhat peculiarly. It starts counting with zero when counting from the front and with one when counting from the back.

But lists are not the only things in Python that can be accessed in the way shown in ???. Python introduces a special class of types the sequence types.

Sequences in Python
Definition 2.5.6. Python has more types that behave just like lists, they are called sequence types.
\triangleright The most important sequence types for IWGS are lists, strings and ranges.
▷ Definition 2.5.7. A range is a finite sequence of numbers it can conveniently be constructed by the range function: range((⟨start⟩), ⟨(stop⟩), ⟨(step⟩)) constructs a range from ⟨(start⟩) (inclusive) to ⟨(stop⟩) (exclusive) with step size ⟨(step⟩).
▷ Example 2.5.8. Lists can be constructed from ranges:
>>> list(range(1,6,2)) [1,3,5]
range(1,6,2) makes a "range" from 1 to 6 with step 2, list makes it a list.
FAU : 56 2025-06-05 EXERCISE

Ranges are useful, because they are easily and flexibly constructed for iteration (up next). You may ask yourselves, why Python has a special data structure for ranges. The main reason is that we can treat them more efficiently than lists. Consider the range constructed by range(1,100000000), i.e. the numbers between 1 and a billion. If we were to represent this as a list, then this would probably take most of the memory available on your laptop, even if we do not do anything with it. But in the range, Python only needs to actually create those parts of the range that it actually needs. Say we want to access the 1000th element, then the Python interpreter can just compute that as 1+1000 when it needs to (and free the memory when that is no longer needed); in particular, the interpreter does even have to create all the intermediate elements.



CHAPTER 2. INTRODUCTION TO PROGRAMMING

▷ Example 2.5.11. Lists a	and strings can also act as sequer	ices.	(try it)
<pre>print("Let me reverse sou x = input("please type so</pre>	mething for you!")		
for i in reversed(list(x)):	Sincetining.)		
FAU	57	2025-06-05	COMMERCIAL DE LA COMMERCIA DE LA COMMERCIA DE LA COMMERCIA DE LA COMPACIÓN DE

But lists are not the only data structure for collections of objects. Python provides others that are organized slightly differently for different applications. We give a particularly useful example here: dictionaries.

Python Dictionaries		
▷ Definition 2.5.12. A dicting where we call k the key and	onary is an unordered colle v the value.	ection of ordered pairs (k,v) ,
In Python dictionaries are w mata, and the value is separately and the value is sep	written with curly brackets, rated from the key by a co	, pairs are separated by com- lon.
▷ Example 2.5.13. Dictiona	ries can be used for variou	s purposes,
painting = {	dict_de_en = dt", "Maus": "r Watch", "Ast": "bra "Klavier": " } can be nested, e.g. for a l	= { mouse", anch", "piano" ist of paintings. enum = { 1: "copy", 2: "paste", 3: "adapt" }
FAU	58	2025-06-05

Dictionaries give "keyed access" to collections of data: we can access a value via its key. In particular, we do not have to remember the position of a value in the collection.

Interacting with Dict	ionaries						
⊳ Example 2.5.14 (Diction)	onary operations).						
▷ painting["title"] retur	rns the value for the key "tit	le" in the dictionary painting.					
⊳ painting["title"]="De original Dutch	painting["title"]="De Nachtwacht" changes the value for the key "title" to its original Dutch (or adds item "title": "De Nachtwacht")						
⊳ Example 2.5.15 (Print	▷ Example 2.5.15 (Printing Keys and Values).						
keys	values	key/value pairs					
for × in thisdict.keys(): print(×) ▷ More dictionary commar	<pre>for x in thisdict.values(): print(x) ds:</pre>	<pre>for x, y in thisdict.items(): print(x, y)</pre>					
▷ if 《key》 in 《dict》 ch▷ painting.pop("title")	ecks whether $\langle\!\langle key angle\!\rangle$ is a key removes the "title" item fro	/ in 《dict》. om painting.					

40

FAU : 59	2025-06-05	
----------	------------	--

Note that thisdict.keys has a short form: we can just iterate over the keys using for x in thisdict:.

2.5.2 Input and Output

The next topic of our stroll through Python is one that is more practically useful than intrinsically interesting: file input/output. Together with the regular expressions this allows us to write programs that transform files.



We now fix some of the nomenclature surrounding files and file systems provided by most operating system. Most programming languages provide their own bindings that allow to manipulate files.



Many operating systems use files as a primary computational metaphor, also treating other resources like files. This leads to an abstraction of files called streams, which encompass files as well as e.g. keyboards, printers, and the screen, which are seen as objects that can be read from (keyboards) and written to (e.g. screens). This practice allows flexible use of programs, e.g. re-directing a the (screen) output of a program to a file by simply changing the output stream.

41

Now we can come to the Python bindings for the file input/output operations. They are rather straightforward.



The only interesting thing is that we have to declare our intentions when we opening a file. This allows the file system to protect the files against unintended actions and also optimize the data transfer to the storage devices involved.

Let us now look at some examples to fortify our intuition about what we can do with files in practice.

```
Disk Input/Output in Python (continued)
 ▷ Example 2.5.21 (Reading a file linewise).
        >>> f.readline()
                                                 >>> for line in f:
        'This is the first line of the file.\n'
                                                 ... print(line, end='')
        >>> f.readline()
                                                 This is the first line of the file.
        'Second line of the file\n'
                                                 Second line of the file
        >>> f.readline()
 \triangleright If you want to read all the lines of a file in a list you can also use list(f) or
   f.readlines().
    For reading a Python file we use the import((basename)) instruction
    ▷ it searches for the file 《basename》.py, loads it, interprets it as Python code,
      and directly executes it.
    ▷ primarily used for loading Python libraries
                                                             (additional functionality)
```

⊳ also useful for load	ling Python-encoded data	(e.g. dicti	onaries)
FAU	63	2025-06-05	CCCCC AND A DECEMBER OF A DECE

The code snippet on the right of Example 2.5.21 show that files can be iterated over using a for loop: the file object is implicitly converted into a sequences of strings via the readline method.

2.5.3 Functions and Libraries in Python

We now come to a general device for organizing and modularizing code provided by most programming languages, including Python. Like variables, functions give names to Python objects – here fragments of code – and thus make them reusable in other contexts.

Functions in Python (Introduction)
Observation: Sometimes programming tasks are repetitive print("Hello Peter, how are you today? How about some IWGS?") print("Hello Roxana, how are you today? How about some IWGS?") print("Hello Frodo, how are you today? How about some IWGS?)
 Idea: We can automate the repetitive part by functions. Example 2.5.22.We encapsultate the greeting functionality in a function:
<pre>def greet (who): print("Hello ",who," how are you today? How about some IWGS?") greet("Peter") greet("Roxana") greet("Frodo") greet(input ("Who are you?"))</pre>
 and use it repeatedly. Functions can be a very powerful tool for structuring and documenting programs (if used correctly)
FAU : 64 2025-06-05
Functions in Python (Example)
Example 2.5.23 (Multilingual Greeting). Given a value for lang def greet (who): if lang == 'en': print("Hello ",who," how are you today? How about some IWGS?") elif lang == 'de': print("Sehr geehrter ",who,", wie geht's heute? Wie waere es mit IWGS?") we can even localize (i.e. adapt to the language specified in lang) the greeting.
FAU : 65 2025-06-05 CONTRACTOR

We can now make the intuitions above formal and give the exact Python syntax of functions.



We now come to a peculiarity of an object-oriented language like Python: it treats types as first-class entities, which can be defined by the user – they are called classes then. We will not go into that here, since we will not need classes in IWGS, but have have to briefly talk about methods, which are essentially functions, but have a special notation.

Python provides two kinds of function-like facilities: regular functions as discussed above and methods, which come with Python classes. We will not attempt a presentation of object oriented programming and its particular implementation in Python this would be beyond the mandate of the IWGS course – but give a brief introduction that is sufficient to use methods.

Functions vs. Methods in Python
▷ There is another mechanism that is similar to functions in Python. (we briefly introduce it here to delineate)
\triangleright Background: Actually, the types from ??? are classes,
▷ Definition 2.5.25. In Python all values belong to a class, which provide special functions we call methods. Values are also called objects, to emphasise class aspects. Method application is written with dot notation: (obj) . $(meth)$ ($(args)$) corresponds to $(meth)$ ((obj) , $(args)$).
▷ Example 2.5.26. Finding the position of a substring
>>> s = 'This is a Python string' # s is an object of class 'str'
>>> type(s) <class 'str'=""> # see, I told you so >>> s index('Puthen') # det notation (index is a string method)</class>
10
FAU : 67 2025-06-05 E



For the purposes of IWGS, it is sufficient to remember that methods are a special kind of functions that employ the dot notation. They are provided by the class of an object.

It is very natural to want to share successful and useful code with others, be it collaborators in a larger project or company, or the respective community at large. Given what we have learned so far this is easy to do: we write up the code in a (collection of) Python files, and make them available for download. Then others can simply load them via the **import** command.

```
Pvthon Libraries
 \triangleright Idea: Functions, classes, and methods are re usable, so why not package them up
   for others to use.
 ▷ Definition 2.5.28. A Python library is a Python file with a collection of functions,
   classes, and methods. It can be imported (i.e. loaded and interpreted as a Python
   program fragment) via the import command.
 \triangleright There are \ge 150.000 libraries for Python
                                                    (\hat{=} packages on http://pypi.org)
    ▷ search for them at http://pypi.org
                                                         (e.g. 815 packages for "music")
    \triangleright install them with pip install (package name)
    \triangleright look at how they were done
                                                           (all have links to source code)
    \triangleright maybe even contribute back (report issues, improve code, ...) (open source)
e
                                           69
                                                                       2025-06-05
```

The Python community is an open source community, therefore many developers organize their code into libraries and license them under open source licenses.

Software repositories like PyPI (the Python Package Index) collect (references to) and make them for the package manager pip, a program that downloads Python libraries and installs them on the local machine where the **import** command can find them.

2.5.4 A Final word on Programming in IWGS

This leaves us with a final word on the way we will handle prgramming in this course: IWGS is not a programming course, and we expect you to pick up Python from the IWGS and web/book resources. So, recall:



Our very quick introduction to Python is intended to present the very basics of programming and get IWGS students off the ground, so that they can start using programs as tools for the humanities and social sciences.

But there is a lot more to the core functionality Python than our very quick introduction showed, and on top of that there is a wealth of specialized packages and libraries for almost all computational and practical needs.

Chapter 3

Numbers, Characters, and Strings

In our basic introduction to programming above we have convinced ourselves that we need some basic objects to compute with, e.g. Boolean values for conditionals, numbers to calculate with, and characters to form strings for input and output. In this chapter we will look at how these are represented in the computer, which in principle can only store binary digits voltage or no noltage on a wire – which we think of as 1 and 0.

In this chapter we look at the representation of the basic data structures of programming languages (numbers and characters) in the computer; Boolean values ("True" and "False") can directly be encoded as binary digits.

Documents as Digital Objects
Question: how do texts get onto the computer?(after all, computers can only do 0/1)
\triangleright Hint: At the most basic level, texts are just sequences of characters.
\triangleright Answer: We have to encode characters as sequences of bits.
⊳ We will go into how:
▷ documents are represented as sequences of characters,
▷ characters are represented as numbers,
ightarrow numbers are represented as bits (0/1).
FAU : 71 2025-06-05 CONTRACTOR

3.1 Representing and Manipulating Numbers

We start with the representation of numbers. There are multiple number systems, as we are interested in the principles only, we restrict ourselves to the natural numbers – all other number systems can be built on top of these. But even there we have choices about representation, which influence the space we need and how we compute with natural numbers.

The first system for number representations is very simple; so simple in fact that it has been discovered and used a long time ago.

Natural Numbers



In addition to manipulating normal objects directly linked to their daily survival, humans also invented the manipulation of place-holders or symbols. A *symbol* represents an object or a set of objects in an abstract way. The earliest examples for symbols are the cave paintings showing iconic silhouettes of animals like the famous ones of Cro-Magnon. The invention of symbols is not only an artistic, pleasurable "waste of time" for humans, but it had tremendous consequences. There is archaeological evidence that in ancient times, namely at least some 8000 to 10000 years ago, humans started to use tally bones for counting. This means that the symbol "bone with marks" was used to represent numbers. The important aspect is that this bone is a symbol that is completely detached from its original down to earth meaning, most likely of being a tool or a waste product from a meal. Instead it stands for a universal concept that can be applied to arbitrary objects. So far so good, let us see how this would be represented on a computer:

3.1. REPRESENTING AND MANIPULATING NUMBERS



The problem with the unary number system is that it uses enormous amounts of space, when writing down large numbers. We obviously need a better representation. The unary natural numbers are very simple and direct, but they are neither space-efficient, nor easy to manipulate. Therefore we will use different ways of representing numbers in practice.

Positional Number Systems ▷ **Problem:** Find a better representation system for natural numbers. \triangleright Idea: Build a clever code on the unary natural numbers, use position information and addition, multiplication, and exponentiation. \triangleright **Definition 3.1.3.** A positional number system \mathcal{N} is a pair $\langle D, \varphi \rangle$ with $\triangleright D$ is a finite set of b digits; b := #(D) is the base or radix of \mathcal{N} . $\triangleright \varphi \colon D \to [0, b-1]$ is bijective. We extend φ to a bijection between sequences d_k, \ldots, d_0 of digits and natural numbers by setting $\varphi(d_k,\ldots,d_0) := \sum_{i=0}^k \varphi(d_i) \cdot b^i$ We say that the digit sequence $n_b := d_k, \ldots, d_0$ is the positional notation of a natural number n, iff $\varphi(d_k, \ldots, d_0) = n$. \triangleright Example 3.1.4. $\langle \{a, b, c\}, \varphi \rangle$ with with $\varphi(a) := 0$, $\varphi(b) := 1$, and $\varphi(c) := 2$ is a positional number system for base three. We have $\varphi(c, a, b) = 2 \cdot 3^2 + 0 \cdot 3^1 + 1 \cdot 3^0 = 18 + 0 + 1 = 19$ FAU e 2025-06-05 74

If we look at the unary number system from a greater distance, we see that we are not using a very important feature of strings here: position. As we only have one letter in our alphabet, we cannot, so we should use a larger alphabet. The main idea behind a positional number system $\mathcal{N} = \langle D_b, \varphi_b \rangle$ is that we encode numbers as strings of digits in D_b , such that the position matters, and to give these encodings a meaning by mapping them into the unary natural numbers via a mapping φ_b .

Commonly Used Positional Number Systems										
\triangleright Definition 3.1.5. The following positional number systems are in common use.										
	name	Set	Dase	uigits	example					
	unary	\mathbb{N}_1	1	0	000001					
	binary	\mathbb{N}_2	2	0,1	01010001112					
	octal	\mathbb{N}_8	8	0,1,,7	63027 ₈					
	decimal	\mathbb{N}_{10}	10	0,1,,9	162098_{10} or 162098					
	hexadecimal	№ 16	16	0,1,,9,A,,F	FF3A12 ₁₆					

Binary digits are also called bits, and a sequence of eight bits an octet.

 ▷ Notation: Attach the base of N to every number from N. (default: decimal)
 ▷ Trick: Group triples or quadruples of binary digits into recognizable chunks (add leading zeros as needed)
 ▷ 1100011010111002 = 01102 00112 01012 11002 = 635C16 ○ 1100011010111002 = 01102 0012 0112 0102 = 615348 ○ F3A₁₆ = F₁₆ 3₁₆ A₁₆ = 1111001110102, 47218 = 48 78 28 18 = 1001110100012 ○ F3A₁₆ = F₁₆ 3₁₆ A₁₆ = 1111001110102, 47218 = 48 78 28 18 = 1001110100012

We have all seen positional number systems: our decimal system is one (for the base 10). Other systems that important for us are the binary system (it is the smallest non degenerate one) and the octal (base 8) and hexadecimal (base 16) systems. These come from the fact that binary numbers are very hard for humans to scan. Therefore it became customary to group three or four digits together and introduce (compound) digits for these groups. The octal system is mostly relevant for historic reasons, the hexadecimal system is in widespread use as syntactic sugar for binary numbers, which form the basis for electronic circuits, since binary digits can be represented physically by voltage/no voltage.



How to get back to Decimal (or any other system)

- \triangleright **Observation:** ??? specifies how we can get from base *b* representations to decimal. We can always go back to the base *b* numbers.
- \triangleright **Observation 3.1.7.** To convert a decimal number n to base b, use successive integer division (division with remainder) by b:

i := n; repeat (record $i \mod b$, $i := i \dim b$) until i = 0.

 \triangleright Example 3.1.8 (Convert 456 to base 8). Result: 710₈



3.2 Characters and their Encodings: ASCII and UniCode

IT systems need to encode characters from our alphabets as bit strings (sequences of binary digits (bits) 0 and 1) for representation in computers. To understand the current state – the unicode standard – we will take a historical perspective. It is important to understand that encoding and decoding of characters is an activity that requires standardization in multi-device settings – be it sending a file to the printer or sending an e-mail to a friend on another continent. Concretely, the recipient wants to use the same character mapping for decoding the sequence of bits as the sender used for encoding them – otherwise the message is garbled.

We observe that we cannot just specify the encoding table in the transmitted document itself, (that information would have to be en/decoded with the other content), so we need to rely document-external external methods like standardization or encoding negotiation at the metalevel. In this section we will focus on the former.

The ASCII code we will introduce here is one of the first standardized and widely used character encodings for a complete alphabet. It is still widely used today. The code tries to strike a balance between being able to encode a large set of characters and the representational capabilities in the time of punch cards (see below).

<u>The A</u>	The ASCII Character Code															
Definition 3.2.1. The American Standard Code for Information Interchange (ASCII) is a character encoding that assigns characters to numbers 0 to 127.																
Code	0	1	$ \cdots 2$	3	4	5	6	7	8	9	$\cdots A$	$\cdots B$	$\cdots C$	$\cdots D$	$\cdots E$	$\cdots F$
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
1	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
$2 \cdot \cdot \cdot$	ш	1		#	\$	%	&	/	(]	*	+	,	—		/
$3 \cdots$	0	1	2	3	4	5	6	7	8	9	1	;	<	=	>	?
4	0	A	В	C	D	E	F	G	H	I	J	K	L	M	N	0
5	P	Q	R	S	Т	U	V	W	X	Y	Z	l	\]		_
6	-	a	b	с	d	e	Í	g	h	1	J	ĸ		m J	n ~	O
 ▷ The first 32 characters are control characters for ASCII devices like printers. ▷ Motivated by punch cards: The character 0 (0000000₂ in binary) carries no information NUL, (used as dividers) Character 127 (= 111111₂) can be used for deleting (overwriting) last value (cannot delete holes) 																
 The ASCII code was standardized in 1963 and is still prevalent in computers today. (but seen as US centric) 																

Punch cards were the preferred medium for long-term storage of programs up to the late 1970s, since they could directly be produced by card punchers and automatically read by computers.



Up to the 1970s, computers were batch machines, where the programmer delivered the program to the operator (a person behind a counter who fed the programs to the computer) and collected the printouts the next morning. Essentially, each punch card represented a single line (80 characters) of program code. Direct interaction with a computer is a relatively young mode of operation. The ASCII code as above has a variety of problems, for instance that the control characters are mostly no longer in use, the code is lacking many characters of languages other than the English language it was developed for, and finally, it only uses seven bits, where an octet (eight bits) is the preferred unit in information technology. Therefore a whole zoo of extensions were introduced, which — due to the fact that there were so many of them — never quite solved the encoding problem.

Problems with ASCII encoding

- ▷ Problem: Many of the control characters are obsolete by now/ (e.g. NUL,BEL, or DEL)
- \triangleright **Problem:** Many European characters are not represented. (e.g. $\dot{e}, \tilde{n}, \ddot{u}, \beta, ...$)
- ▷ European ASCII Variants: Exchange less-used characters for national ones.
- ▷ Example 3.2.4 (German ASCII). Remap e.g. [→ Ä,]→ Ü in German ASCII ("Apple][" comes out as "Apple ÜÄ")
- ▷ Problem: No cursive Arabic, Asian, African, Old Icelandic Runes, Math,...

⊳ Idea: scalable	Do somethir architecture, s	ng totally eparate	different to	include	all the	world's sc	ripts:	For a
⊳ wha ⊳ a ma	t characters are apping from bit	e available t strings to	, and o characters			ch) (characte	naract er enc	er set) oding)
Fau	:		80			2025-06-05		e

The goal of the UniCode standard is to cover all the worlds scripts (past, present, and future) and provide efficient encodings for them. The only scripts in regular use that are currently excluded are fictional scripts like the elvish scripts from the Lord of the Rings or Klingon scripts from the Star Trek series.

An important idea behind UniCode is to separate concerns between standardizing the character set — i.e. the set of encodable characters and the encoding itself.

Unicode and the Universal Character Set			
Definition 3.2.6 (Twin Standards). A scalable architecture for representing all the worlds writing systems:			
▷ The universal character set (UCS) defined by the ISO/IEC 10646 International Standard, is a standard set of characters upon which many character encodings are based.			
The unicode standard defines a set of standard character encodings, rules for normalization, decomposition, collation, rendering and bidirectional display or- der.			
Definition 3.2.7. Each UCS character is identified by an unambiguous name and an natural number called its code point.			
\triangleright The UCS has 1.1 million code points and nearly 100 000 characters.			
▷ Definition 3.2.8. Most (non-Chinese) characters have code points in [1,65536]: the basic multilingual plane (BMP).			
Definition 3.2.9 (Notation). For code points in the (BMP), four hexadecimal digits are used, e.g. U + 0058 for the character LATINCAPITALLETTERX;			
FAU : 81 2025-06-05 CONTRACTOR			

Note that there is indeed an issue with space-efficient character encodings here. UniCode reserves space for 2^{32} (more than a million) characters to be able to handle future scripts. But just simply using 32 bits for every UniCode character would be extremely wasteful: UniCode-encoded versions of ASCII files would be four times as large.

Therefore UniCode allows multiple character encodings. UTF-32 is a simple 32-bit code that directly uses the code points in binary form. UTF-8 is optimized for western languages and coincides with the ASCII where they overlap. As a consequence, ASCII encoded texts can be decoded in UTF-8 without changes — but in the UTF-8 encoding, we can also address all other unicode characters (using multi-byte characters).

Character Encodings in Unicode



Note how the fixed bit prefixes in the UTF-8 encoding are engineered to determine which of the four cases apply, so that UTF-8 encoded documents can be safely decoded.



XKCD's Take on Recent Unicode Extensions (cont.)



3.3 More on Computing with Strings

We now extend our repertoire on handling and formatting strings in Python: we will introduce string literals, which allow writing complex strings.

```
Playing with Strings and Characters in Python
 ▷ Definition 3.3.1. Python strings are sequences of UniCode characters.
 \triangleright A In Python, characters are just strings of length 1.
 ▷ ord gives the UCS code point of the character, chr character for a number.
 ▷ Example 3.3.2 (Playing with Characters).
   def lc(c) :
       return chr(ord(c) + 32)
   def uc(c) :
       return chr(ord(c) - 32)
   >>> uc('d')
   'D'
   >>> lc('D')
 \triangleright Strings can be accessed by ranges [i:j]
                                                                          ([i] \cong [i:i])
 Example 3.3.3. Taking strings apart and re-assembling them.
   def cap(s) :
       if s == "":
           return "" # base case
       else:
           return uc(s[0]) + cap(s[1:len(s)])
   >>> cap('iwgs')
   'IWGS'
FAU
                                                                               85
                                                                   2025-06-05
```

Example 3.3.3 may be difficult to understand at first. It is a programming technique called recursion, i.e. functions that call themselves from within their body to solve problems by utilizing solutions to smaller instances of the same problem. Recursion can lead to very concise code, but requires some getting-used-to.

In Example 3.3.3 we define a function cap that given a string s returns a string that is constructed by combining the first character uppercased by the uc function with the result of calling the cap function on the rest string – s without the first character. The base case for the recursion is the empty string, where uc also returns the empty string. So let us see what happens in our test cap('iwgs'):

 $cap('iwgs') \rightsquigarrow uc('i')+cap('wgs') \rightsquigarrow 'l'+uc('w')+cap('gs') \rightsquigarrow 'l'+'W'+uc('g')+cap('s') \rightsquigarrow 'lW'+'G'+cap('s') \rightsquigarrow 'lWG'+uc('s')+cap('') \rightsquigarrow 'lWG'+'S'+cap('') \rightsquigarrow 'lWGS'+'' \rightsquigarrow 'lWGS'$

Example 3.3.2 and Example 3.3.3 (or any other examples in this lecture) are not production code, but didactically motivated – to show you what you can do with the objects we are presenting in Python.

In particular, if we "lowercase" a character that is already lowercase – e.g. by lc('c'), then we get out of the range of the UCS code: the answer is x83, which is the character with the hexadecimal code 83 (decimal 131), i.e. the character No Break Here.

In production code (as used e.g. in the Python lower method), we would have some range checks, etc.

String Literals in Python

- ▷ **Problem:** How to write strings including special characters?
- ▷ Definition 3.3.4. A literal is a notation for representing a fixed value for a data structure in source code.
- ▷ Definition 3.3.5. Python uses string literals, i.e character sequences surrounded by one, two, or three sets of matched single or double quotes for string input. The content can contain escape sequences, i.e. the escape character backslash followed by a code character for problematic characters:

Seq	Meaning	Seq	Meaning
//	Backslash (\)	\'	Single quote (')
/"	Double quote (")	\a	Bell (BEL)
\b	Backspace (BS)	\f	Form-feed (FF)
\n	Linefeed (LF)	\r	Carriage Return (CR)
\t	Horizontal Tab (TAB)	\v	Vertical Tab (VT)

In triple-quoted string literals, unescaped newlines and quotes are honored, except that three unescaped quotes in a row terminate the literal.

FAU

2025-06-05

Raw String Literals in Python

- Definition 3.3.6. Prefixing a string literal with a r or R turns it into a raw string literal, in which backslashes have no special meaning.
- ▷ **Note:** Using the backslash as an escape character forces us to escape it as well.
- \triangleright Example 3.3.7. The string "a\nb\nc" has length five and three lines, but the

3.3. MORE ON COMPUTING WITH STRINGS

string r"a\nb\nc" only has length seven and only one line.				
Fau	:	87	2025-06-05	CC) String addition (Sector)

Now that we understand the "theory" of encodings, let us work out how to program with them in Python:

Programming with UniCode strings is particularly simple, strings in Python are UTF-8-encoded UniCode strings and all operations on them are UniCode-based¹. This makes the introduction to UniCode in Python very short, we only have to know how to produce non-ASCII characters, i.e. the characters that are not on regular keyboards.

If we know the code point, this is very simple: we just use UniCode escape sequences.

Unicode in Python				
▷ Remark 3.3.8. The Python string data type is UniCode, encoded as UTF-8. ▷ How to write UniCode characters?: there are five ways				
▷ write them in your editor	(make sure that it uses UTF-8)			
⊳ otherwise use Python escape sequence	es (try it!)			
>>> "\xa3" # Using 8—bit hex valu '\u00A3' >>> "\u00A3" # Using a 16—bit he '\u00A3' >>> "\U000000A3" # Using a 32— '\u00A3' >>> "\N{Pound Sign}" # characte '\u00A3'	ue ex value bit hex value r name			
FAU	2025-06-05			

Note that the discussion about entry methods for unicode characters applies to the bare Python interpreter, not Python-specific text editor modes or user interfaces, which are often helpful by automatically replacing the input by the respective glyphs themselves.

String literals are convenient for creating simple string objects. For more complex ones, we usually want to build them from pieces, usually using the values of variables or the results of functions. This is what f strings are for in Python; we will cover that now.

Formatted String Literals (aka. f-strings)			
Problem: In a program we often want to build strings from pieces that we already have lying around interspersed by other strings.			
Solution: Use string concatenation: >>> course="IWGS"			
>>> students=6*11 >>> "The " + course + " course has " + <mark>str</mark> (studen [.] 'The IWGS course has 66 students'	ts) + " students"		
▷ We can do better! (mixing blank	s and quotes is error-prone)		

¹Older programming languages have ASCII strings only, and UniCode strings are supplied by external libraries.



FAU

89

2025-06-05



3.4 More on Functions in Python

We now extend our repertoire of dealing with functions in Python.

In a sense, we now know all we have to about Python function: we can define them and apply them to arguments. But Python offers us much more: Python

- treats functions as "first-class objects", i.e. entities that can be given to other functions as arguments, and can be returned as results.
- provides more ways of passing arguments to a function than the rather rigid way we have seen above. This can be very convenient and make code more readable.

We will cover these features now. The main motivation for this is that they are widely used in programming and being able to read them is important for collaborating with experienced programmers and reading existing code.

We digress to the internals of functions that make them even more powerful. It turns out that we do not have to give a function a name at all.

Anonymous Functions (lambda)

- ▷ **Observation 3.4.1.** A Python function definition combines making a function object with giving it a name.
- ▷ Definition 3.4.2. Python also allows to make anonymous functions via the function

3.4. MORE ON FUNCTIONS IN PYTHON

literal lambda for	function objects:			
lambda p_1,\ldots,p_n :	: «(expr)			
⊳ Example 3.4.3.	The following tw	o Python fragments are	e equivalent:	
	def cube (x): x*x*x	cube = lambda x: x*	«X*X	
The right one is just a variable assignment that assigns a function object to the variable cube. (In fact Python uses the right one internally)				
Question: Why use anonymous functions?				
Answer: We may not want to invent (i.e. waste) a name if the function is only used once. (examples on the next slide)				
FAU		91	2025-06-05	

Anonymous functions do not seem like a big deal at first, but having a way to construct a function that can be used in any expression, is very powerful as we will see now.

Higher-Order Functions in Python ▷ **Definition 3.4.4.** We call a function a higher order function, iff it takes a function as argument. ▷ **Definition 3.4.5.** map and filter are built-in higher order functions in Python. They take a function and a list as arguments. ightarrow map(f,L) returns the list of f-values of the elements of L. \triangleright filter(p,L) returns the sub-list L' of those l in L, such that p(l)=True. ▷ Example 3.4.6. Mapping over and filtering a list >>> li = [5, 7, 22, 97, 54, 62, 77, 23, 73, 61] >>> list(map(lambda x: x*2 , li)) [10, 14, 44, 194, 108, 124, 154, 46, 146, 122] >>> list(filter(lambda x: (x%2 != 0), li)) [5, 7, 97, 77, 23, 73, 61] FAU 2025-06-05 92

Admittedly, in our example, we could also have defined a named function twice and then mapped that over li:

def twice (x): x*x map twice li

But the code from Example 3.4.6 is more compact. Once we get used to the programming idiom and understand it, it becomes quite readable.

Another important feature of Python functions is flexible argument passing. This allows to define functions that supply complex behaviors – for which we need to set many parameters but simple calling patterns – which is good to hide complexity from the programmer.

The first argument passing feature we want to discuss is the use of keyword arguments, which gets around the problem of having to remember the position of an argument of a multi-parameter function.

Argument Passing in Python: Keyword Arguments			
▷ Definition 3.4.7. The last $k \le n$ of n parameters of a function can be keyword arguments of the form $p_i = \langle val \rangle_i$: If no argument a_i is given in the function call, the default value $\langle val \rangle_i$ is taken.			
▷ Example 3.4.8. The head of the open function is			
def open(file, mode='r', buffering=-1, encoding=None, errors=None, newline=None, closefd=True, opener=None)			
Even if we only call it with open("foo"), we can use parameters like mode or opener in the body; they have the corresponding default value.			
We can also give more arguments via keywords, even out of order			
open("foo", buffering=1, mode="+a")			
FAU : 93 2025-06-05 CONTRACTOR			

BTW: The opener argument of open is a function, and often an anonymous function is used if it is specified.

The next feature is dual to the last: instead of letting the caller leave out some arguments, we allow the caller more, which is then bound to a list parameter.

Actually the star operator can be used in other situations as well, consider for instance

>>> numbers = [2, 1, 3, 4, 7] >>> more_numbers = [*numbers, 11, 18]

3.4. MORE ON FUNCTIONS IN PYTHON

```
>>> print(*more_numbers, sep=', ')
2, 1, 3, 4, 7, 11, 18
```

Here we have used the star operator twice: First to pass the list numbers as arguments to the list constructor and a second time to pass the extended list more_numbers to the print function. Finally, we can combine the ideas from the last two to make keyword arguments flexary.

```
Argument Passing in Python: Flexible Keyword Arguments
 ▷ Definition 3.4.13. Python functions can take keyword arguments:
   if k is a sequence of key/value pairs then def f(p_1, \dots, p_n, **k) binds the keys to
   values in the body of f.
 ▷ Example 3.4.14.
   def kw args(farg, **kwargs):
       print (f"formal arg: {farg}")
       for key in kwargs :
           print (f"another keyword arg: {key}: {kwargs[key]}")
   >>> kw args(1, myarg2="two", myarg3=3)
   formal arg: 1
   another keyword arg: myarg2 : two
   another keyword arg: myarg3 : 3
FAU
                                                                          95
                                                               2025-06-05
```

Just as for the flexible arity case above, we have an operator that unpacks argument structures, here a dictionary.



Disclaimer: The last couple of features of Python functions are a bit more advanced than would usually be expected from a Python programming introduction in a course such as IWGS. But one of the goals of IWGS is to empower students to be able to read Python code of more

experienced authors. And that kind of code may very well contain these features, so we need to cover them in IWGS.

So the last couple of slides should be considered as an "early exposure for understanding" rather than "essential to know for IWGS" content.

3.5 Regular Expressions: Patterns in Strings

Now we can come to the main topic of this section: regular expressions, A domain-specific language for describing string patterns. Regular expressions are extremely useful, but also quite cryptical at first. They should be understood as a powerful tool, that relies on a language with a very limited vocabulary. It is more important to understand what this tool can do and how it works in principle than memorizing the vocabulary – that can be looked up on demand.

Problem: Text/Data File Manipulation				
Problem 1 (Information Extraction): We often want to extract information from large document collections, e.g.				
 ▷ e-mail addresses or dates from collected correspondencesrtts ▷ dates and places from newsfeeds ▷ links from web pages 				
Problem 2 (Data Cleaning): The representation in data files is often too noisy and inconsistent for directly importing into an application; e.g.				
 ▷ standardizing different spellings of e.g. city names, (Nuremberg vs. Nürnberg) ▷ eliminating higher UniCode characters, when the application only accepts ASCII, ▷ separating structured texts into data blocks. (e.g. in x-separated lists) 				
\triangleright Enabling Technology: Specifying text/data fragments \rightsquigarrow regular expressions.				
FAU : 97 2025-06-05				

There are several dialects of regular expression languages that differ in details, but share the general setup and syntax. Here we introduce the Python variant and recommend [PyRegex] for a cheat-sheet on Python regular expressions (and an integrated regular expression tester).

 Regular Expressions, see [Pyt]

 ▷ Definition 3.5.1. A regular expression (also called regular expression) is a formal expression that specifies a set of strings.

 ▷ Definition 3.5.2 (Meta-Characters for Regexps).

	char	denotes]	
	•	any single character (except a newline)]	
	^	beginning of a string]	
	\$	end of a string]	
	$[\ldots]/[^{\ldots}]$	any single character in/not in the brackets]	
	$[x-y]/[^x-y]$	any single character in/not in range x to y		
	[]	marks a capture group]	
	$\setminus n$	the n^{th} captured group		
		disjunction]	
	*	matches preceding element zero or more times]	
	+	matches preceding element one or more times]	
	?	matches preceding element zero or one times		
	$\{n,m\}$	matches the preceding element between n and m times		
	\S/\s	non-/whitespace character]	
	\W/\w	non-/word character]	
	\D/\d	non-/digit (not only 0-9, but also e.g. arabic digits)]	
			-	
All other characters match themselves to match e.g. a? escape with a $\langle \cdot \rangle$				
Fau	:	98 2025-06-05		
			=	

Let us now fortify our intuition with some (simple) examples and a more complex one.



As we have seen regular expressions can become quite cryptic and long (cf. e.g. ???), so we need help in developing them. One way is to use one of the many regexp testers online

Playing with Regular Expressions

 \triangleright If you want to play with regular expressions, go e.g. to http://regex101.com
e		Jacobs University × 🙀 http://www.c?id=program × 👧 Online regex tester and d ×	+	2	
(🗲 📀 regex101.com/#pyth	on 🗸 🖓 🔞 🕻	Google	• • • • • • • • • • • • • • • • • • •	
	📄 Services 👻 📄 News 👻 🚺	MathWeb - 🚺 Jacobs - 🛄 AG - 🛄 Lists - 🛄 TRAC - 🛄 Rotary - 🛄 TNTBase - 🛄 Haus -			
	regular expressions 101	>_ regex tester 📥 community 🙊 irc		regex101 💲 donate 🛛 contact 🔺 bug reports & suggestions	
	FLAVOR pcre (php) javascript python TOOLS format regex (req	REGULAR DARRESSION " "Elblot "geixsu TEST STRING the rat bit the KOE	О	EXPLANATION	
	 code generator regex debugger post to community VERSION CONTROL Save regex ACCOUNT ⊕ log in 			MATCH INFORMATION No match groups were extracted. This means that your pattern matches but there were no kapturing groups in it that matched anything in the subject string.	
ł	SETTINGS display whitespace way long lines colorize syntax use dark theme use minimal view			QUICK REFERENCE MOST USED TOKENS FULL REFERENCE MOST USED TOKENS mast used tokens A single character of a, b or c (a) Itokens A character except a, b or c (b) CATEGODES A character except a, b or c (b) general tokens A character not in the ranges := [set]	
EAII		SUBSTITUTION	0	d anchors A character in the range: a [e3A-2]	0
	:	100		2025-06-05	IMERICHTS RESE

After covering regular expressions in the abstract, we will see how they are integrated into programming languages to solve problems. Of course we take Python as an example.

Regular Expressions in Python
We can use regular expressions directly in Python by importing the re module (just add import re at the beginning)
\triangleright As Python has UniCode strings, regular expressions support UniCode as well.
▷ Useful Python functions that use regular expressions.
▷ re.findall(《pat》,《str》): Return a list of non-overlapping matches of 《pat》 in 《str》.
>>> re.findall(r"[h c r]at",'the cat ate the rat on the mat') ['cat','rat']
$ ightarrow$ re.sub($\langle\!\langle pat \rangle\!\rangle$, $\langle\!\langle sub \rangle\!\rangle$, $\langle\!\langle str \rangle\!\rangle$): Replace substrings that match $\langle\!\langle pat \rangle\!\rangle$ in $\langle\!\langle str \rangle\!\rangle$ by $\langle\!\langle sub \rangle\!\rangle$.
>>> re.sub(r'\sAND and\s', '', 'Baked Beans and Spam')'Baked Beans Spam'
ightarrow re.split(((pat)), ((str)): Split ((str)) into substrings that match ((pat)).
>>> re.split(r'\s+','When shall we three meet again?'))
['When', 'shall', 'we', 'three', 'meet', 'again?'] >>> re.split(r'\s+ \? \. ! , : ; ', 'When shall we three meet again?')) ['When', 'shall', 'we', 'three', 'meet', 'again']
FAU : 101 2025-06-05 EXECT

As regular expressions form a special language for describing sets of strings, it is not surprising that they are used in all kinds of searching, splitting, and substring replacement operations. As the language of regular expressions is well standardized, these more or less work the same in all programming languages, so what you learn for Python, you can re-use in other languages.

We will now see what we can do with regular expressions in a practical example. You should consider it as a "code reading/understanding" exercise, not think of it as something you should (easily) be able to do yourself. But ??? could serve as a quarry of ideas for things you can do to texts with regular expressions.



This program is just a series of stepwise regular expression computations that are assigned to the variable s. For the last one, we use the **lambda** operator that constructs a function as an argument (the second) to re.sub. We use the anonymous functions because this function is only used once. This worked well, so we just continue along these lines.

Example: Correcting and Anonymizing Documents (cont.)
▷ Example 3.5.6 (Document Cleanup (continued)).
\triangleright next we make abbreviations for regular expressions to save space
c = [A-Z]
= a-z
\triangleright remove capital letters in the middle of words
$s = re.sub(f''({I})({c}+)({I})'',$
s) #
\triangleright and we cross-out for official public versions of government documents,
$s = re.sub(f''({c}{l} + ({c}{l} * (.?))?{c}{l} +)", #$
s)



We show the whole program again, to see that it is relatively small (thanks to the very compact – if cryptic – regular expressions), when we leave out all the comments.



Chapter 4

Documents as Digital Objects

In this chapter we take a first look at documents and how they are represented on the computer.

4.1 Representing & Manipulating Documents on a Computer

Now that we can represent characters as bit sequences, we can represent text documents. In principle text documents are just sequences of characters; they can be represented by just concatenating them.

Electronic Documents

- Definition 4.1.1. An electronic document is any media content that is intended to be used via a document renderer, i.e. a program or computing device that transforms it into a form that can be directly perceived by the end user.
- **Example 4.1.2.** PDFs, digital images, videos, audio recordings, web pages, ...
- Definition 4.1.3. An electronic document that contains a digital encoding of textual material that can be read by the end user by simply presenting the encoded characters is called digital text.
- ▷ **Definition 4.1.4.** Digital text is subdivided into plain text, where all characters carry the textual information and formatted text, which also contains instructions to the document renderer.

▷ Example 4.1.5. Python programs are plain text, PDFs are formatted.

FAU : 105 2025-06-05

We will now establish a nomenclature for giving instructions to a document renderer. This has originated from movable (lead) type based typesetting but carries over well to electronic documents.

Document Markup

Definition 4.1.6. Document markup (or just markup) is the process of adding control words (special character sequences also called markup codes) to a plain text



There are many systems for document markup, ranging from informal ones as in Example 4.1.7 that specify the intended document appearance to humans – in this case the printer – to technical ones which can be understood by machines but serving the same purpose.

Markup is by no means limited to visual markup for documents intended for printing as Example 4.1.7 may suggest. There are aural markup formats that instruct document renderers that transform documents to audio streams of e.g. reading speeds, intonation, and stress.

We now come to another aspect of electronic documents: We mostly interact with them in the form of files. Again, we fix our nomenclature.

File Types

- ▷ Observation 4.1.10. We mostly encounter electronic documents in the form of files on some storage medium.
- ▷ Definition 4.1.11. A text file is a file that contains text data, a binary file one that contains binary data

Definition 4.1.12. Text files are often processed as a sequence of text lines (or just lines), i.e. sub string separated by the line feed character U + 000A; LINEFEED(LF). The line number is just the position in the sequence.

- ▷ Remark 4.1.13. Text files are usually encoded with ASCII, ISO Latin, or increasingly UniCode encodings like UTF-8.
- **Example 4.1.14.** Python programs are stored in text files.
- \triangleright In practice, text files are often processed as a sequence of text lines (or just lines), i.e. sub strings separated by the line feed character U + 000A; LINEFEED(LF). The line number is just the position in the sequence.



Remark 4.1.15. Plain text is different from formatted text, which includes markup code, and binary files in which some portions must be interpreted as binary data (encoded integers, real numbers, digital images, etc.)

As we have seen above, it does not take much to render a text file: we only need to guess the right encoding scheme so we can decode the file and show the character sequence to the user. Indeed the UNIX cat just prints the contents of a text file to a shell. But we need much more, we need tools with which we can compose and edit text files; we do this with text editors, which we will discuss now.

Text Editors

- Definition 4.1.16. A text editor is a program used for rendering and manipulating text files.
- **Example 4.1.17.** Popular text editors include
 - ▷ Notepad is a simple editor distributed with MSWindows.
 - emacs and vi are powerful editors originating from UNIX and optimized for programming.
 - ▷ sublime is a sophisticated programming editor for multiple operating systems.
 - ▷ EtherPad is a browser-based real-time collaborative editor.
- ▷ **Example 4.1.18.** Even though it can save documents as text files, MSWord is not usually considered a text editor, since it is optimized towards formatted text; such "editors" are called word processors.

FAU : 108 2025-06-05

What text editors do for text files, word processors do for other electronic documents.



FAU : 109 2025-06-0	
---------------------	--

Before we go on, let us first get into some basics: how do we measure information, and how does this relate to units of information we know.

4.2 Measuring Sizes of Documents/Units of Information

Having represented documents as sequences of characters, we can use that to measure the sizes of documents. In this section we will have a look at the underlying units of information and try to get an intuition about what we can store in files.

We will take a very generous stance towards what a document is, in particular, we will include pictures, audio files, spreadsheets, computer aided designs,

Units for Information			
Observation: The smallest unit of information is knowing the state of a sys with only two states.	stem		
▷ Definition 4.2.1. A bit (a contraction of "binary digit") is the basic unit of capacity of a data storage device or communication channel. The capacity of a system which can exist in only two states, is one bit (written as 1b)			
▷ Note: In the ASCII encoding, one character is encoded as 8b, so we introduce another basic unit:			
\triangleright Definition 4.2.2. The byte is a derived unit for information capacity: $1B = 8$	3b.		
FAU : 110 2025-06-05	O		

From the basic units of information, we can make prefixed units for prefixed units for larger chunks of information. But note that the usual SI unit prefixes are inconvenient for application to information measures, since powers of two are much more natural to realize.

Larger Units of Information via Binary Prefixes

- ▷ We will see that memory comes naturally in powers to 2, as we address memory cell by binary number, therefore the derived information units are prefixed by special prefixes that are based on powers of 2.
- ▷ Definition 4.2.3 (Binary Prefixes). The following binary unit prefixes are used for information units because they are similar to the SI unit prefixes.

prefix	symbol	2^n	decimal	~SI prefix	Symbol
kibi	Ki	2^{10}	1024	kilo	k
mebi	Mi	2^{20}	1048576	mega	Μ
gibi	Gi	2^{30}	1.074×10^{9}	giga	G
tebi	Ti	2^{40}	1.1×10^{12}	tera	Т
pebi	Pi	2^{50}	1.125×10^{15}	peta	Р
exbi	Ei	2^{60}	1.153×10^{18}	exa	Е
zebi	Zi	2^{70}	1.181×10^{21}	zetta	Ζ
yobi	Yi	2^{80}	1.209×10^{24}	yotta	Υ

4.2. MEASURING SIZES OF DOCUMENTS/UNITS OF INFORMATION

Note: The correspondence works better on the smaller prefixes; for yobi vs. yotta there is a 20% difference in magnitude.
 The SI unit prefixes (and their operators) are often used instead of the correct binary ones defined here.
 Example 4.2.4. You can buy hard-disks that say that their capacity is "one terabyte", but they actually have a capacity of one tebibyte.

Fau

111

CC Same rights reserved

2025-06-05

Let us now look at some information quantities and their real-world counterparts to get an intuition for the information content.

<u>How n</u>	<u>nuch Informatic</u>	n?	
	Bit (b)	binary digit 0/1	
	Byte (B)	8 bit	
	2 Bytes	A UniCode character in UTF.	
	10 Bytes	your name.	
	Kilobyte (kB)	1,000 bytes OR 10^3 bytes	
	2 Kilobytes	A Typewritten page.	
	100 Kilobytes	A low-resolution photograph.	
	Megabyte (MB)	1,000,000 bytes OR 10^6 bytes	
	1 Megabyte	A small novel or a 3.5 inch floppy disk.	
	2 Megabytes	A high-resolution photograph.	
	5 Megabytes	The complete works of Shakespeare.	
	10 Megabytes	A minute of high-fidelity sound.	
	100 Megabytes 1 meter of shelved books.		
	500 Megabytes A CD-ROM.		
	Gigabyte (GB)	1,000,000,000 bytes or 10^9 bytes	
	1 Gigabyte	a pickup truck filled with books.	
	20 Gigabytes	A good collection of the works of Beethover	ı.
	100 Gigabytes	A library floor of academic journals.	
Fau	:	112 2025-06-05	SAWE HIR HERE WED

How much Information?

71

Γ	Terabyte (TB)	1,000,000,000,000 bytes or 10^{12} bytes	
	1 Terabyte	50000 trees made into paper and printed.	
	2 Terabytes	An academic research library.	
	10 Terabytes	The print collections of the U.S. Library of Congress.	
	400 Terabytes	National Climate Data Center (NOAA) database.	
	Petabyte (PB)	1,000,000,000,000,000 bytes or 10^{15} bytes	
	1 Petabyte	3 years of EOS data (2001).	
	2 Petabytes	All U.S. academic research libraries.	
	20 Petabytes	Production of hard-disk drives in 1995.	
	200 Petabytes	All printed material (ever).	
	Exabyte (EB)	1,000,000,000,000,000,000 bytes or 10^{18} bytes	
	2 Exabytes	Total volume of information generated in 1999.	
5 Exabytes		All words ever spoken by human beings ever.	
	300 Exabytes	All data stored digitally in 2007.	
	Zettabyte (ZB)	1,000,000,000,000,000,000,000 bytes or 10^{21} bytes	
	2 Zettabytes	Total volume digital data transmitted in 2011	
	100 Zettabytes	Data equivalent to the human Genome in one body.	
-			
Fai	y .	113 2025-06-05	CC Somerication exercise

The information in this table is compiled from various studies, most recently [HL11].

Note: Information content of real-world artifacts can be assessed differently, depending on the view. Consider for instance a text typewritten on a single page. According to our definition, this has ca. 2kB, but if we fax it, the image of the page has 2MB or more, and a recording of a text read out loud is ca. 50MB. Whether this is a terrible waste of bandwidth depends on the application. On a fax, we can use the shape of the signature for identification (here we actually care more about the shape of the ink mark than the letters it encodes) or can see the shape of a coffee stain. In the audio recording we can hear the inflections and sentence melodies to gain an impression on the emotions that come with text.

4.3 Hypertext Markup Language

WWW documents have a specialized document type that mixes markup for document structure with layout markup, hyper-references, and interaction. The HTML markup elements always concern text fragments, they can be nested but may not otherwise overlap. This essentially turns a text into a document tree. In IWGS, we discuss HTML mostly as a way to build interfaces of web applications. Therefore we will prioritize those aspects of HTML that have to do with "programming documents" over the creation of nice-looking web pages. Therefore we will pick up the notion of nested text fragments marked up by well-bracketed tags and elements in section 4.4 and generalize these ideas to XML as a general representation paradigm for semi-structured data in ???.

We will also postpone the discussion of cascading stylesheets, which have evolved as the dominant technology for the specification of presentation (layout, colors, and fonts) for marked-up documents, to ???.

4.3.1 Introduction

HTML was created in 1990 and standardized in version 4 in 1997 [RHJ98]. Since then the WWW has evolved considerably from a web of static web pages to a Web in which highly dynamic web pages become user interfaces for web-based applications and even mobile applets. HTML5 standardized the necessary infrastructure in 2014 [Hic+14].



The thing to understand here is that HTML uses the characters \langle , \rangle , and / to delimit the markup. All markup is in the form of tags, so anything that is not between \langle and \rangle is the textual content.

We will not give a complete introduction to the various tags and elements of the HTML language here, but refer the reader to the HTML recommendation [Hic+14] and the plethora of excellent web tutorials. Instead we will introduce the concepts of HTML markup by way of examples.

The best way to understand HTML is via an example. Here we have prepared a simple file that shows off some of the basic functionality of HTML.

A very first HTML Example (Source)
<html xmlns="http://www.w3.org/1999/xhtml"></html>
<head></head>
<title>A first HTML Web Page</title>
<body></body>
<h1>Anatomy of a HTML Web Page</h1>
<h3>Michael Kohlhase FAU Erlangen Nuernberg</h3>
<h2 id="intro">1. Introduction</h2>
This is really easy, just start writing.
<h2>3. Main Part: show off features</h2>
We can can markup text styles inline.
And we can make itemizations:

vith a list item
and another one

<h2>4. Conclusion<th></th><th></th></h2>		
$\langle \mathbf{p} \rangle$ As we have seen in	the introducti	on this
was very easy.		
FAU :	115	2025-06-05

The thing to understand here is that HTML markup is itself a well-balanced structure of begin and end tags. That wrap other balanced HTML structures and – eventually – textual content. The HTML recommendation [Hic+14] specifies the visual appearance expectation and interactions afforded by the respective tags, which HTML-aware software systems – e.g. a web browser – then execute. In the next slide we see how FireFox displays the HTML document from the previous.

A very first HTM	L Example (Result)		
	 ▲ first HTML Web Page × + ● ← → C i file:///Users/kohlhas ···· ⊙ ☆ » = Anatomy of a HTML Web Page Michael Kohlhase FAU Erlangen Nürnberg 1. Introduction This is really easy, just start writing. 3. Main Part: show off features We can can markup text styles inline. And we can make itemizations: with a list item and another one 3. Conclusion As we have seen in the introduction this was very easy. 		
FAU	116	2025-06-05	CC State Rights Reserved

4.3.2 Interacting with HTML in Web Broswers

In the last slide, we have seen FireFox as a document renderer for HTML. We will now introduce this class of programs in general and point out a few others.

Web Browsers
▷ Definition 4.3.4. A web browser is a software application for retrieving (via HTTP), presenting, and traversing information resources on the WWW, enabling users to view web pages and to jump from one page to another.
Definition 4.3.5. A web browser usually supplies user tools like
▷ history that gives the user access to web pages visited earlier and

4.3. HYPERTEXT MARKUP LANGUAGE

▷ bookmark to remember web pages.					
Definition 4.3.6. A web browser usually supplies developer tools like					
 ▷ the console that logs system-level events in the browser and ▷ an inspector that gives access to the structure and content of the DOM. 					
> Practical Browser Tools:					
 ▷ Status Bar: security info, page load progress ▷ Favorites (bookmarks) ▷ View Source: view the code of a web page ▷ Tools/Internet Options, history, temporary Internet files, home page, auto complete, security settings, programs, etc. 					
Example 4.3.7 (Common Browsers).					
 MSInternetExplorer is an once dominant, now obsolete browser for MSWindows. Edge is provided by Microsoft for MSWindows.(replaces MSInternetExplorer) FireFox is an open source browser for all platforms, it is known for its standards compliance. Safari is provided by Apple for macOS and MSWindows. Chrome is a lean and mean browser provided by Google Inc. (very common) WebKit is a library that forms the open source basis for Safari and Chrome. 					
EAU : 117 2025-06-05 EXAMPLE					

Let us now look at a couple of more advanced tools available in most web browsers for dealing with the underlying HTML document.

Browser Tools for dealing with HTML, e.g. in FireFox

 \vartriangleright Hit Control-U to see the page source in the browser



We have used **FireFox** as an example here, but these tools are available in some form in all major browsers the browser vendors want to make their offerings attractive to web developers, so that web pages and web applications get tested and debugged in them and therefore work as expected.

4.3.3 A Worked Example: The Contact Form

After this simple example, we will come to a more complex one: a little "contact form" as we find on many web sites that can be used for sending a message to the owner of the site. Let us only look a the design of the form document before we go into the interaction facilities afforded it.

4.3. HYPERTEXT MARKUP LANGUAGE

HTML in Practice: Worked Ex	ample
▷ Make a design and "paper prototype"	of the page:
Rease type your e-	in a Mersæge mæi/actores:
 ▷ Put the intended text into a file: cont Contact Please enter a message: Your e—mail address: xx @ xx.de Send message ▷ Load into your browser to check the s 	act.html:
/Users/kohlhase/localmh/	MathHub// × +
$(\epsilon) \rightarrow C $	i file:///Users/kohlhase/localmh/MathHub
Getting Started 📄 FAU 📄 Services	🗋 News 🗎 MathWeb 🗎 AG 📄 Rotary 🏠 N
Contact Please type in a message: Yo	ur e-mail address: xx @ xx.de Send message
ho Add title, paragraph and button mark	up:
< title >Contact <b title> < h2 >Please enter a message: <b h2> < h3 >Your e—mail address: xx @ xx. < button >Send message <b button>	Contact ← → C ⓓ ⓒ Getting Started ☐ FAU ☐ Services ☐ Please type in a message: Your contact eMail address: xx @ xx.de Send message
ightarrow Add input fields and breaks:	



After designing the functional (what are the text blocks) structure of the contact form, we will need to understand the interaction with the contact form.



We current ignore the form data (the part after the ?)

> We will come to the full story of processing actions later.

2025-06-05

Unfortunately, we can only see what the browser sends to the server at the current state of play, not what the server does with the information. But we will get to this when we take up the example again.

For the moment, we made use of the fact that we can just specify the page contact—after.html, which the browser displays next. That ignores the query part and – via a form tags of its own gets the user back to the original contact form.

More useful types of Input fields	
▷ Radio buttons: type="radio" (grouped	by name attribute)
<input name="gender" type="radio" value="male"/> Male <input name="gender" type="radio" value="female"/> Female <input name="gender" type="radio" value="other"/> Other	○ Male ○ Female ○ Other
▷ Check boxes: type="checkbox"	
My major is <input name="major" type="checkbox" value="cs"/> Computer Science <input name="major" type="checkbox" value="dh"/> Digital Humanities <input name="major" type="checkbox" value="other"/> Other	S
My major is Computer Science Digital Humanities Other	
▷ File selector dialogs (interaction is system specific here for	or MacOS Mojave)
Upload your resume <input name="resume" type="file"/>	
Upload your resume Browse No file selected.	
\triangleright Drop down menus: select and option	
Which animal do you like? <select name="animals"> </select>	
	2025-06-05

4.4 Documents as Trees

We have concentrated on HTML as a document type for interactive multimedia documents. Before we progress, we want to discuss an important feature: all practical document types that employ control words are in some sense well-bracketed. Well-bracketed structures are well-understood in CS and mathematics: they are called trees and come with a rich and useful collection of descriptive concepts and tools. We will present the concepts in this section and the tools they enable in ???.





4.4. DOCUMENTS AS TREES

Trees are well understood mathematical objects and tree data structures are very commonly used in computer science and programming. As such they have a well-developed nomenclature, which we will introduce now.





Why are trees written upside-down?: The main answer is that we want to draw tree diagrams in text. And we naturally start drawing a tree at the root. So, if a tree grows from the root and we do not exactly know the tree height, then we do not know how much space to leave. When we write trees upside down, we can directly start from the root and grow the tree downward as long as we need. We will keep to this tradition in the IWGS course.



We will now make use of the tree structure for computation. Even if the computing tasks we pursue here may seem a bit abstract, they show very nicely how tree algorithms typically work.



82

def height (tree):	def maxh (I):
return $maxh(tree[1:]) + 1$	if I == []:
	return 0
height([1,[2,[[4],[5]]],[3,[[6],[7]]]])	else
>>> 3	<pre>return max(height(l[0]),maxh(l[1:]))</pre>
FAU	126 2025-06-05 CONTRACTOR

Let use have a closer look at Example 4.4.9. The algorithm consists of two functions:

- 1. height, which computes the height of an input tree by delegating the computation of the maximal height of its children to maxh and then incrementing the value by 1.
- 2. maxh, which takes a list of trees and computes the maximum of their heights by calling height on the first input tree and then comparing with the maximal height of the remaining trees.

Note that maxh and height each call the other. We call such functions mutually recursive. Here this behavior poses no problem, since the arguments in the recursive calls are smaller than the inputs: for maxh it is the rest list, and for height the "list of children" of the input tree.

Example 4.4.9 was complex for two reasons: mutual recursion and the somewhat cryptic encoding of trees as lists of lists of integers. We claim that recursive programming is "not a bug, but a feature", as it allows to succinctly capture the "divide-and-conquer" approach afforded by trees. For the cryptic encoding of trees we can do better.



Again, we have two mutually recursive functions: weight that takes a tree, and wsum that takes a list and the recursion goes analogously. Only that this time, the list of children is a dictionary value and the calls are clearer. The only real difference, is that in wsum we have to add up the weight of the head of the list an the joint sum of the rest list.



storing marked up electronic documents as trees together with a standardized set of access methods for manipulating them.
 Idea: When a web browser loads a HTML page, it directly parses it into a DOM and then works exclusively on that. In particular, the HTML document is immediately discarded; documents are rendered from the DOM.

4.5 An Overview over XML Technologies

We have seen that many of the technologies that deal with marked-up documents utilize the tree-like structure of (the DOM) of HTML documents. Indeed, it is possible to abstract from the concrete vocabulary of HTML that implements the intended layout of hypertexts and the function of its fragments, and build a generic framework for document trees. This is what we will study in this section.

4.5.1 Introduction to XML

XML (EXtensible Markup Language)					
Definition 4.5.1. XML (short for Extensible Markup Language) is a framework for markup formats for documents and structured data.					
▷ Tree representation language		(begin/er	d brackets)		
\triangleright Restrict instances by <i>Doc. Ty</i>	pe Def. (DTD) or .	Schema	(Grammar)		
▷ Presentation markup by style	files	(XSL: <mark>X</mark> ML <mark>S</mark> tyle	Language)		
\triangleright Intuition: XML is extensible H	TML				
▷ logic annotation (<i>markup</i>) instead of presentation!					
ightarrow many tools available: parsers, compression, data bases,					
▷ conceptually: transfer of trees instead of strings.					
\triangleright details at http://w3c.org	(XML is standardi	ize by the WWW (Consortium)		
FAU .	129	2025-06-05	CONTRACTOR OF THE SECOND		

The idea of XML being an "extensible" markup language may be a bit of a misnomer. It is made "extensible" by giving language designers ways of specifying their own vocabularies. As such XML does not have a vocabulary of its own, so we could have also it an "empty" markup language that can be filled with a vocabulary.

$\frac{\text{XML} \text{ is Everywhere (E.g. We}}{\text{WL}}$	eb Pages)
Example 4.5.2. Open web page f you get the following text:	ile in FireFox, then click on $View \searrow PageSource$ (showing only a small part and reformatting)
<pre><html xmlns="http://www.w3.org/1999</td><td>/xhtml"></html></pre>	
<pre><title>Michael Kohlhase</title> <meta <="" name="generator" pre=""/></pre>	



Now we see an example of an XML file that is used for communicating data in a machine-readable, but human-understandable way.

XML is Everywhere (E.g. Catalogs)
Example 4.5.4 (The NYC Galleries Catalog). A public XML file at https://data.cityofnewyork.us/download/kcrmj9hh/application/xml
x<b ml version="1.0" encoding="UTF-8"?> <museums> <museum></museum></museums>
<name>American Folk Art Museum</name>
<pre><phone>212-265-1040</phone></pre>
<address>45 W. 53rd St. (at Fifth Ave.)</address>
<closing>Closed: Monday</closing>
<rates>admission: \$9; seniors/students, \$7; under 12, free</rates>
<specials></specials>
Pay—what—you—wish: Friday after 5:30pm;
refreshments and music available
<museum></museum>
<name>American Museum of Natural History</name>
<pre><phone>212-769-5200</phone></pre>
<address>Central Park West (at W. 79th St.)</address>
<closing>Closed: Thanksgiving Day and Christmas Day</closing>
FAU : 131 2025-06-05 EXECUTIVE

This XML uses an ad hoc markup language: Every <museum> element represents one museum in New York City (NYC). Its children convey the detailed information as "key value pairs". And now, if you still need proof that XML is really used almost everywhere, here is the ultimate example.



4. Other files have packaging information, images, and other objects.
▲ This is huge and offensively ugly.
> But you have everything you wanted and more
> In particular, you can process the contents via a program now.





86

4.5.2 Computing with XML in Python

We have claimed above that the tree nature of XML documents is one of the main advantages. Let us now see how Python makes good on this promise.

We use the external lxml library [LXMLa] in IWGS, even though the Python distribution includes the standard library ElementTree library [ET] for dealing with XML. lxml subsumes ElementTree and extends it by functionality for XPath and can parse a large set of HTML documents even though they are not valid XML. This makes lxml a better basis for practical applications in the Digital Humanities.

Acknowledgements: Many of the examples and the flow of exposition in the next slides has been adapted from the lxml tutorial [LXMLc].







This method of "manually" producing XML trees in memory by applying etree methods may seem very clumsy and tedious. But the power of lxml lies in the fact that these can be embedded in Python programs. And as always, programming gives us the power to do things very efficiently.



But XML documents that only have elements, are boring; let's do XML attributes next. Recall that attributes are essentially string-valued key/value pairs. So what could be more natural than treating them like dictionaries.



Recall that we could use Python dictionaries for iterating over in a for loop. We can do the same for attributes:

Computing with XML in Python (Attributes; continued) \triangleright We can access attributes by the keys, values, and items methods, known from dictionaries: >>> sorted(root.keys()) ['hello', 'interesting'] >>> for name, value in sorted(root.items()): ... print(f'{name} = {value}') hello = 'Huhu' interesting = 'totally' \triangleright \triangle To get a 'real' dictionary, use the attrib method (e.g. to pass around) >>> attributes = root.attrib Note that attributes participates in any changes to root and vice versa. \triangleright \triangle To get an independent snapshot of the attributes that does not depend on the XML tree, copy it into a dict: >>> d = dict(root.attrib)>>> sorted(d.items()) [('hello', 'Guten Tag'), ('interesting', 'totally')] FAU e 2025-06-05 139

The last two items touch a somewhat delicate subject in programming. Mutable an immutable data structures: the former can be changed in place as we have above with the .set method, and the latter cannot. Both have their justification and respective advantages. Immutable data structures are "safe" in the sense that they cannot be changed unexpectedly by another part of the program, they have the disadvantage that every time we want to have a variant, we have to copy the whole object. Mutable ones do not – we can change in place – but we have to be very careful about who accesses them when.

This is also the reason why we spoke of "dictionary-like interface" to XML trees in lxml: dictionaries are immutable, while XML trees are not.

The main remaining functionality in XML is the treatment of text. XML treats text as special kinds of node in the tree: text nodes. They can be treated just like any other node in the XML tree in the etree library.

Computing with XML in Python (Text nodes)

> XML elements can contain text: we use the text property to access and set it.
>>> root = etree.Element("root")
>>> root.text = "TEXT"
>>> print(root.text)
TEXT
>>> etree.tostring(root)
b'croot>TEXT</root>'

far: we programmatically construct an HTML tree.

Case Study: Creating an HTML document
▷ We create nested html and body elements
>>> html = etree.Element("html") >>> body = etree.SubElement(html, "body")
\triangleright Then we inject a text node into the latter using the .text property.
>>> body.text = "TEXT"
\triangleright Let's check the result
>>> etree.tostring(html) b' <html><body>TEXT</body></html> '
\triangleright We add another element: a line break and check the result
>>> br = etree.SubElement(body, "br")
>>> etree.tostring(html)
D < Itimi>< body>1 Ex1< br/>>
arappi Finally, we can add trailing text via the .tail property
>> br.tail = "TAIL"
>>> etree.tostring(html) b' <html><body>TEXT TAIL</body></html> '
FAU : 141 2025-06-05 CONTRACTOR

Note the use of the .tail property here? While the .text property can be used to set "all" the text in an XML element, we have to use the .tail property to add trailing text (e.g. after the $\langle br \rangle >$ element).

Notwithstanding the "Python power" argument from above, there are situations, where we just want to write down XML fragments and insert them into (programmatically created) XML trees. lxml as functionality for this: XML literals, which we introduce now.

Computing with XML in Python (XML Literals)
Definition 4.5.8 We call any string that is well-formed XML an XML literal
\triangleright We can use the XML function to read XML literals.
>>> root = etree.XML(" <root>data</root> ")
The result is a first-class element tree, which we can use as above
>>> print(root.tag) root
>>> etree.tostring(root)
b' <root>data</root> '
BTW, the fromstring function does the same.
▷ There is a variant html that also supplies the necessary HTML decoration.

90



4.5.3 XML Namespaces

We now come to a topic that is considered very difficult, confusing, and un-necessary by many people: XML namespaces. But it really is not, if you approach it with an open mind. Indeed it is probably what you would have come up with if you had been presented with the problem of mixing vocabularies, which is in turn a consequence of the fact that XML is used pervasively in the computing world and especially in Digital Humanities, where we often need to aggregate semi-structured data from multiple sources (and this multiple XML vocabularies).



This is an excerpt from the document metadata which AcrobatDistiller saves along with each PDF document it creates. It contains various kinds of information about the creator of the document, its title, the software version used in creating it and much more. Document metadata is useful for libraries, bookselling companies, all kind of text databases, book search engines, and generally all institutions or persons or programs that wish to get an overview of some set of books,

documents, texts. The important thing about this document metadata text is that it is not written in an arbitrary, PDF proprietary format. Document metadata only make sense if these metadata are independent of the specific format of the text. The metadata that MSWord saves with each Word document should be in the same format as the metadata that Amazon saves with each of its book records, and again the same that the British library uses, etc.

We will now reflect what we have seen in Example 4.5.9 and fully define the namespacing mechanisms involved. Note that these definitions are technically involved, but conceptually quite natural. As a consequence they should be read more with an eye towards "what are we trying to achieve" than the technical details.



4.5.4 XPath: Specifying XML Subtrees

One of the great advantages of viewing marked-up documents as trees is that we can describe subsets of its nodes.



4.5. AN OVERVIEW OVER XML TECHNOLOGIES



An XPath processor is an application or library that reads an XML file into a DOM and given an XPath expression returns (pointers to) the set of nodes in the DOM that satisfy the expression.



To see that XPath is not just a plaything, we will now look at at a typical example where we can identify useful subtrees in a large HTML document: the Wikipedia page on paintings by Leonardo da Vinci.

 XPath Example: Scraping Wikipedia

 ▷ Example 4.5.15 (Extracting Information from HTML).



If the task of writing an XPath for extracting the 50+ titles from this page does not convince you as worth learning XPath for, consider that Wikipedia has ca. 30 such lists, which apparently have exactly the same tree structure, so the XPath developed once for da Vinci, probably works for all the others as well.

Chapter 5

Web Applications

In this chapter we will see how we can turn HTML pages into web-based applications that can be used without having to install additional software.

For that we discuss the basics of the World Wide Web as the client server architecture that enables such applications. Then we take up the contact form example to get an understanding how information is passed between client and server in interactive web pages. This motivates a discussion of server-side computation of web pages that can react to such information. A discussion of CSS styling shows how to make the web pages that are generated can be made visually appealing. We conclude the chapter by a discussion of client-side computation that allows making web pages interactive without recurring to the server. **Excursion:** The World Wide Web as we introduce it here is based on the Internet infrastructure and protocols. In some places it may be useful to read up on this insection A.1.

5.1 Web Applications: The Idea



We have seen that web applications are a common way of building application software. To understand how this works let us now have a look at the components.



To understand web applications, we will first need to understand

- 1. how we can express web pages in HTML and (see ???) interact with them for data input (we recap this in ???),
- 2. the basics of how the World Wide Web works as a distribution framework (see ???),
- 3. how we can generate HTML documents programmatically (in our case in Python; see ???) as answer pages, and finally
- 4. how we can make HTML pages dynamic by client side manipulation (see ???).

5.2 Basic Concepts of the World Wide Web

We will now present a very brief introduction into the concepts, mechanisms, and technologies that underlie the World Wide Web and thus web applications, which are our interest here.

5.2.1 Preliminaries

The WWW is the hypertext/multimedia part of the internet. It is implemented as a service on top of the internet (at the application level) based on specific protocols and markup formats for documents.

The Internet and the Web

▷ Definition 5.2.1. The Internet is a global computer network that connects hundreds of thousands of smaller networks.

- Definition 5.2.2. The World Wide Web (WWW) is an open source information space where electronic documents and other web resources are identified by URLs, interlinked by hypertext links, and can be accessed via the Internet.
- ▷ Intuition: The WWW is the multimedia part of the internet, they form critical infrastructure for modern society and commerce.
- ▷ The internet/WWW is huge:

Year	Web	Deep Web	eMail
1999	21 TB	100 TB	11TB
2003	167 TB	92 PB	447 PB
2010	????	?????	?????

▷ We want to understand how it works. (services and scalability issues)
Image: 150 2025-06-05

Given this recap we can now introduce some vocabulary to help us discuss the phenomena.

Concepts of the World Wide Web > Definition 5.2.3. A web page is a document on the WWW that can include multimedia data and hyperlinks. ▷ **Note:** Web pages are usually marked up in in HTML. ▷ Definition 5.2.4. A web site is a collection of related web pages usually designed or controlled by the same individual or organization. \triangleright A web site generally shares a common domain name. ▷ **Definition 5.2.5.** A hyperlink is a reference to data that can immediately be followed by the user or that is followed automatically by a user agent. ▷ **Definition 5.2.6.** A collection text documents with hyperlinks that point to text fragments within the collection is called a hypertext. The action of following hyperlinks in a hypertext is called browsing or navigating the hypertext. \triangleright In this sense, the WWW is a multimedia hypertext. Fau e 2025-06-05 151

5.2.2 Addressing on the World Wide Web

The essential idea is that the World Wide Web consists of a set of resources (documents, images, movies, etc.) that are connected by links (like a spider-web). In the WWW, the links consist of pointers to addresses of resources. To realize them, we only need addresses of resources (much as we have IP numbers as addresses to hosts on the internet).

Uniform Resource Identifier (URI), Plumbing of the Web

> Definition 5.2.7. A uniform resource identifier (URI) is a global identifiers of local

or network-retrievable docum uniform syntax (grammar) de	ents, or media files fined in RFC-3986 [(<mark>web resources</mark>). BLFM05].	URIs adhere a
A URI is made up of the follo	wing components:		
▷ a scheme that specifies the specifies	e protocol governin	g the resource,	
▷ an authority: the host (au	thentication there)	that provides the i	resource,
▷ a path in the hierarchicall	y organized resource	es on the host,	
⊳ a query in the non-hierarc	hically organized pa	rt of the host data	, and
▷ a fragment identifier in th	e resource.		
<pre>Example 5.2.8. The followin http://example.com:8</pre>	ng are two example 8042/over/ther	URIs and their cor e?name=ferre	nponent parts: t#nose
\/ \	/\	/ \	/ \/
scheme authority	path	query f	ragment
,	·l_、		
/ \ / mailto:michael.kohlł	\ nase@fau.de		
Note: URIs only identify de (e.g. in a browser).	ocuments, they do n	ot have to provide	access to them
FAU	152	2025-0	06-05

The definition above only specifies the structure of a URI and its functional parts. It is designed to cover and unify a lot of existing addressing schemes, including URLs (which we cover next), ISBN numbers (book identifiers), and mail addresses.

In many situations URIs still have to be entered by hand, so they can become quite unwieldy. Therefore there is a way to abbreviate them.

Relative URIs						
⊳ Defin in fron	Definition 5.2.9. URIs can be abbreviated to relative URIs; missing parts are filled in from the context.					
⊳ Exam	▷ Example 5.2.10. Relative URIs are more convenient to write					
	relative URI	abbreviates	in context]		
	#foo	$\langle current - file \rangle #foo$	curent file			
	bar.txt	file:///home/kohlhase/foo/bar.txt	file system	1		
	/bar/bar.html	http://example.org/bar/bar.html	on the web	1		
▷ Definition 5.2.11. To distinguish them from relative URIs, we call URIs absolute URIs.						
FAU	:	153	2025-06-05			

The important concept to grasp for relative URIs is that the missing parts can be reconstructed from the context they are found in: the document itself and how it was retrieved.

For the file system example, we are assuming that the document is a file foo.html that was loaded from the file system – under the file system URI file:///home/kohlhase/foo/foo.html – and for the web example via the URI //example.org/foo/foo.html. Note that in the last example, the relative URI ../bar/ goes up one segment of the path component (that is the meaning of ../), and specifies the file bar.html in the directory bar.

But relative URIs have another advantage over absolute URIs: they make a web page or web site easier to move. If a web site only has links using relative URIs internally, then those do not mention e.g. authority (this is recovered from context and therefore variable), so we can freely move the web-site e.g. between domains.

Note that some forms of URIs can be used for actually locating (or accessing) the identified resources, e.g. for retrieval, if the resource is a document or sending to, if the resource is a mailbox. Such URIs are called "uniform resource *locators*", all others "uniform resource *locators*".

Uniform Resource Names and Locators	
▷ Definition 5.2.12. A uniform resource locator (URL) is a URI that gives access to a web resource, by specifying an access method or location. All other URIs are called uniform resource name (URN).	
▷ Idea: A URN defines the identity of a resource, a URL provides a method for finding it.	
▷ Example 5.2.13. The following URI is a URL http://kwarc.info/kohlhase/index.html	(try it in your browser)
▷ Example 5.2.14. urn:isbn:978-3-540-37897-6 only identifies [Koh06] (it is in the library)	
URNs can be turned into URLs via a catalog service, e.g. http://wm-urn.org/ urn:isbn:978-3-540-37897-6	
Note: URIs are one of the core feature considered to be the plumbing of the WWW	es of the web infrastructure, they are /. (direct the flow of data)
FAU : 154	2025-06-05

Historically, started out as URLs as short strings used for locating documents on the internet. The generalization to identifiers (and the addition of URNs) as a concept only came about when the concepts evolved and the application layer of the internet grew and needed more structure. Note that there are two ways in URI can fail to be resource locators: first, the scheme does not support direct access (as the ISBN scheme in our example), or the scheme specifies an access method, but address does not point to an actual resource that could be accessed. Of course, the problem of "dangling links" occurs everywhere we have addressing (and change), and so we will neglect it from our discussion. In practice, the URL/URN distinction is mainly driven by the scheme part of a URI, which specifies the access/identification scheme.




5.2.3 Running the World Wide Web

The infrastructure of the WWW relies on a client-server architecture, where the servers (called web servers) provide documents and the clients (usually web browsers) present the documents to the (human) users. Clients and servers communicate via the HTTPs and HTTPSs protocols. We give an overview via a concrete example before we go into details.



The web browser communicates with the web server through a specialized protocol, the hypertext transfer protocol, which we cover now.

HTTP: Hypertext Transfer Protocol
Definition 5.2.21 The University Transfer Durth rel (UTTD) is an explication
> Deminition 5.2.21. The Hypertext Transfer Protocol (HTTP) is an application
layer protocol for distributed, collaborative, hypermedia information systems.

100

5.2. BASIC CONCEPTS OF THE WORLD WIDE WEB

\triangleright June 1999: HTTP/1.1 is defined in RFC 2616 [Fie+99].				
Preview/Recap: HTTP is used by a client (called user agent) to access web web resources (addressed by uniform resource locators (URLs)) via a HTTP request. The web server answers by supplying the web resource (and metadata).				
Definition 5.2 prominent)	2.22. Most importar	nt HTTP request methods.	(5 more less	
GET	Requests a represent	tation of the specified resource.	safe	
PUT	Uploads a represent	ation of the specified resource.	idempotent	
DELETE	Deletes the specified resource.			
POST	Submits data to be form) to the identifi	e processed (e.g., from a web ed resource.		
 Definition 5.2.23. We call a HTTP request safe, iff it does not change the state in the web server. (except for server logs, counters,; no side effects) Definition 5.2.24. We call a HTTP request idempotent, iff executing it twice has the same effect as executing it once. 				
ightarrow HTTP is a sta	teless protocol.	(very memory efficie	nt for the server.)	
FAU		157 20	025-06-05	

Finally, we come to the last component, the web server, which is responsible for providing the web page requested by the user.

Web Servers > Definition 5.2.25. A web server is a network program (a server in a client server architecture of the WWW) that delivers web resources to and receives content from clients via the Hypertext Transfer Protocol (HTTP). ▷ Example 5.2.26 (Common Web Servers). \triangleright apache is an open source web server that serves about 50% of the WWW. ▷ nginx is a lightweight open source web server. (ca. 35%) ▷ IIS is a proprietary web server provided by Microsoft Inc. ▷ Definition 5.2.27. A web server can host – i.e serve web resources for multiple domains (via configurable hostnames) that can be addressed in the authority components of URLs. This usually includes the special hostname localhost which is interpreted as "this computer". ▷ Even though web servers are very complex software systems, they come preinstalled on most UNIX systems and can be downloaded for MSWindows [Xam]. Fau 158 2025-06-05

Now that we have seen all the components we fortify our intuition of what actually goes down the net by tracing the HTTP messages.



5.3 Recap: HTML Forms Data Transmission

The first two requirement for web applications above are already met by HTML in terms of HTML forms (see slide 120 ff.). Let us recap and extend²

Recap HTML Forms: Submitting Data to the Web Server			
Recall: HTML forms collect data via named input elements, the submit event triggers a HTTP request to the URL specified in the action attribute.			
> Example 5.3.1. Forms contain input fields and explanations.			
<form action="login.html" method="get" name="input"> Username: <input name="user" type="text"/> Password: <input name="pass" type="password"/> <input type="submit" value="Submit"/></form>			
yields the following in a web browser:			
Username: Password: Submit			

EdN:2

5.3. RECAP: HTML FORMS DATA TRANSMISSION



We can now use the tools any modern browser supplies to check up on this claim. In fact, using the browser tools is essential for advanced web development. Here we use the web console, that monitors any activity, to check upon what really happens when we interact with the web page.

Checking	up on the Transmission	
⊳ Let's vei	rify the claims above using browser tools	(here the web console)
⊳ Loading	the file and filling in the form:	(console logs file URI)
	← → C ² ① file:///Users/kohlhase// ···· ♡ ☆ III\ Image: Started Image: make a glossary stu FAU Everyices News MathWeb	 Image: Second se
	Username: mkohihase Password: Submit	
	X Q-Find in page	f 12 matches
	There is a start of the start of t	CSS XHR Requests
	»	<u>(</u>)
⊳ After su	bmitting the form: (console	logs the HTTP request)



A side effect of re-playing our development in the browser is that we see another type of input field: A password field, which hides user input from un-authorized eyes. We also see that the GET request incorporates the form data which contains the password into the URI of the request, which is visible to everyone on the web. We will come back to this problem later.

Let us now look at the data transmission mechanism in more detail to see what is actually transmitted and how.

HTML Forms and Form Data Transmission ▷ We specify the HTTP communication of HTML forms in detail. ▷ **Definition 5.3.2.** The HTML form element groups the layout and input elements: \triangleright <form action="(URI)" method="(reg)"> specifies the form action in terms of a HTTP request $\langle req \rangle$ to the URI $\langle URI \rangle$. \triangleright The form data consists of a string $\langle data \rangle$ of the form $n_1 = v_1 \& \cdots \& n_k = v_k$, where $> n_i$ are the values of the name attributes of the input fields \triangleright and v_i are their values at the time of submission. ▷ <input type="submit" .../> triggers the form action: it composes a HTTP request \triangleright If $\langle (req) \rangle$ is get (the default), then the browser issues a GET request $\langle (URI) \rangle \langle (data) \rangle$. \triangleright If $\langle (req) \rangle$ is post, then the browser issues a POST request to $\langle (URI) \rangle$ with document content (data). \triangleright We now also understand the form action, but should we use GET or POST. FAU 162 2025-06-05

To understand whether we should use the GET or POST methods, we have to look into the details, which we will now summarize.

Practical Differences between HTTP GET and POST					
▷ Using GET vs. POST in HTML Forms:					
		GET	POST		
	Caching	possible	never		
	Browser History	Yes	never		
	Bookmarking	Yes	No		
	Change Server Data	No	Yes		
	Size Restrictions	$\leq 2KB$	No		
	Encryption	No	HTTPS		
 Upshot: HTTP GET is more convenient, but less potent. Always use POST for sensitive data! (passwords, personal data, etc.) GET data is part of the URI and thus unencrypted, POST data via HTTPS is. 					
FAU	163			2025-06-05	

5.4 Generating HTML on the Server

As the WWW is based on a client server architecture, computation in web applications can be executed either on the client (the web browser) or the server (the web server). For both we have a special technology; we start with computation on the web server.

Server-Side Scripting: Programming Web pages
ightarrow Idea: Why write HTML pages if we can also program them! (easy to do)
Definition 5.4.1. A server-side scripting framework is a web server extension that generates web pages upon HTTP requests.
Example 5.4.2. perl is a scripting language with good string manipulation facilities. PERL CGI is an early server-side scripting framework based on this.
Example 5.4.3. Python is a scripting language with good string manipulation facilities. And bottle WSGI is a simple but powerful server-side scripting framework based on this.
Observation: Server-side scripting frameworks allow to make use of external resources (e.g. databases or data feeds) and computational services during web page generation.
▷ Observation: A server-side scripting framework solves two problems:
 making the development of functionality that generates HTML pages convenient and efficient, usually via a template engine, and
2. binding such functionality to URLs the routes, we call this routing.
FAU : 164 2025-06-05 TATAL

We will look at the second problem: routing first. There is a dedicated Python library for that.

5.4.1 Routing and Argument Passing in Bottle

We wil now introduce the bottle library, which supplies a lightweight web server and server-side scripting framework implemented in Python. It is already installed on the JuptyerLab cloud IDE at http://jupyter.kwarc.info. To install it on your laptop, just type pip install bottle in a shell.

The Web Server and Routing in Bottle WSGI
Definition 5.4.4. Serverside routing (or simply routing) is the process by which a web server connects a HTTP request to a function (called the route function) that provides a web resource. A single URI path/route function pair is called a route.
\triangleright The bottle WSGI library supplies a simple Python web server and routing.
 The run(《keys》) function starts the web server with the configuration in 《keys》. The @route decorator connects path components to Python function that return strings. Decorators change functions. A decorator @route(《path》) augments the following function f to answer to HTTP requests to the 《path》 and return f's return value.
▷ Example 5.4.5 (A Hello World route) for localhost on port 8080
from bottle import route, run
@route('/hello') def hello(): return "Hello IWGS!"
run(host='localhost', port=8080, debug=True)
This web server answers to HTTP GET requests for the URL http://localhost: 8080/hello
FAU : 165 2025-06-05 CONTRACTOR

Let us understand Example 5.4.5 line by line: The first line imports the library. The second establishes a route with the name hello and binds it to the Python function hello in line 3 and 4. The last line configures the bottle web server: it serves content via the HTTP protocol for localhost on port 8080.

So, if we run the program from Example 5.4.5, then we obtain a web server that will answer HTTP GET requests to the URL http://localhost:8080/hello with a HTTP answer with the content Hello IWGS!.

To keep the example simple, we have only returned a text string; A realistic application would have generated a full HTML page (see below).

In the last line of Example 5.4.5, we have also configured the **bottle** web server to use "debug mode", which is very helpful during early development.

In this mode, the **bottle** web server is much more verbose and provides helpful debugging information whenever an error occurs. It also disables some optimisations that might get in your way and adds some checks that warn you about possible misconfiguration.

Note that debug mode should be disabled in a production server for efficiency.

But we can do more with routes!

Dynamic Routes in Bottle

▷ Definition 5.4.6. A dynamic route is a route annotation that contains named wildcards, which can be picked up in the route function. \triangleright Example 5.4.7. Multiple @route annotations per route function f are allowed \rightsquigarrow the web application uses f to answer multiple URLs. @route('/') @route('/hello/<name>') **def** greet(name='Stranger'): return (f'Hello {name}, how are you?') With the wildcard <name> we can bind the route function greet to all paths and via its argument name and customize the greeting. **Concretely:** A HTTP GET request to ▷ http://localhost is answered with Hello Stranger, how are you?. > http://localhost/hello/MiKo is answered with Hello MiKo, how are you?. Requests to e.g http://localhost/hello or http://localhost/hello/prof/ kohlhase lead to errors. (404: not found) EAU 0 166 2025-06-05

Often we want to have more control over the routes. We can get that by filters, which can involve data types and/or regular expressions.



We have already seen above that we want to use HTTP GET and POST request for different facets of transmitting HTML form data to the web server. This is supported by bottle WSGI in two ways: we can specify the HTTP method of a route and we have access to the form data (and other aspects of the request).



Recall that we have already seen most of this in slide 160. The only new thing is that we return the HTML as a string in the route function as a request to a HTTP GET request. Now comes the interesting part: the form uses the POST method in the form action and we have to specify a route for that. Recall from ??? that this allows for encrypted transmission, so we are less naive than our solution from slide 160.

```
Bottle Request: Dealing with POST Data
 ▷ Recall: from a HTML form we get a GET or POST request with form data
   n_1 = v_1 \& \cdots \& n_k = v_k
                                (here user=mkohlhase&login=noneofyourbusiness)
 ▷ Bottle WSGI provides the request object for dealing with HTTP request data.
 \triangleright Example 5.4.13 (Login 2).
                                  Continuing from Example 5.4.12: we parse the
   request transmitted request and check password information:
   @post('/login') # or @route('/login', method='POST')
   def do login():
       username = request.forms.get('username')
      password = request.forms.get('password')
      if check login(username, password):
           return "Your login information was correct."
      else:
          return "Login failed."
```

We assume a Python function check login that checks authentication credential and authenticator, and keeps a list of logged in users.

5.4. GENERATING HTML ON THE SERVER



The main new thing in Example 5.4.13 is that we use the request.forms.get method to query the request object that comes with the HTTP request triggering the route for the form data.

5.4.2 Templating in Python via STPL

In IWGS, we use Python for programming, so let us see how we would generate HTML pages in Python.

What would we do in Python
<pre>b Example 5.4.14 (HTML Hello World in Python). print("<html>") print("<body>Hello world</body>") print("</html>")</pre>
▷ Problem 1: Most web page content is static (page head, text blocks, etc.)
$\label{eq:product} \begin{split} & \vdash \textbf{Example 5.4.15 (Python Solution).} \ \dots \text{ use Python functions:} \\ & \textbf{def htmlpage (t,b):} \\ & f'' < html> < head> < title> \{t\} < / title> < / head> < body> \{b\} < / body> < / html> '' \\ & htmlpage("Hello", "Hello IWGS") \end{split}$
▷ Problem 2: If HTML markup dominates, want to use a HTML editor (mode),
▷ e.g. for HTML syntax highlighting/indentation/completion/checking
▷ Idea: Embed program snippets into HTML. (only execute these, copy rest)
FAU : 170 2025-06-05 CONTRACTOR

We will now formalize and toolify the idea of "embedding code into HTML". What comes out of this idea is called "templating". It exists in many forms, and in most programming languages.

Template Processing for HTML				
\triangleright Definition 5.4.16. A template engine (or template processor) for a document format <i>F</i> is a program that transforms templates, i.e. strings or files (a template file) ith a mixture of program constructs and <i>F</i> markup, into a <i>F</i> strings or <i>F</i> documents by executing the program constructs in the template (template processing).				
Note: No program code is left in the resulting web page after generation. (important security concern)				
▷ Remark: We will be most interested in HTML template engines.				
Observation: We can turn a template engine into a server-side scripting framework by employing the URIs of template files on a server as routes and extending the web server by template processing.				
Example 5.4.17. PHP (originally "Programmable Home Page Tools") is a very successful server-side scripting framework following this model.				
FAU : 171 2025-06-05 CONTRACTOR				

Naturally, Python comes with a template engine in fact multiple ones. We will use the one from the bottle web application framework for IWGS.

$\operatorname{stpl:}$ the "Simple Templa	ite Engine'' fr	om Bottle	
▷ Definition 5.4.18 . Bottle W: plate Engine) that processes th (documentation at [STPL])	SGI supplies the t ne STPL (Simple ⁻	emplate engine stpl (Simple Template Language) format.	Tem-
▷ Definition 5.4.19. A templat transforms templates, i.e. strin and F markup, into F-strings of in the template (template proce	e engine for a doc ngs or files throug or <i>F</i> -documents by ressing).	ument format F is a program h a mixture of program cons r executing the program cons	n that structs structs
 stpl uses the template funct program objects into a templat 	ion for template e; it returns a for	processing and {{}} to a matted unicode string.	embed
>>> template('Hello {{name] u'Hello World!'	}}!', name='World	d')	
>>> my_dict={'number': '12 >>> template('I live at {{nur u'I live at 123 Fake St., Fakevi	23', 'street': 'Fake nber}} {{street}} lle'	St.', 'city': 'Fakeville'} , {{city}}', **my_dict)	
FAU	172	2025-06-05	

The stpl template function is a powerful enabling basic functionality in Python, but it does not satisfy our goal of writing "HTML with embedded Python". Fortunately, that can easily be built on top of the template functionality:



next: a line of pythor<br % course = "Informatische Some plain text in betv <% # A block of python code course = name.title().strip %> More plain text	verkzeuge" veen % f	for item in basket: {{item}} end >	
FAU	173	2025-06-05	CC SCALE OF STATE OF STATE

So now, we have template files. But experience shows that template files can be quite redundant; in fact, the better designed the web site we want to to create, the more fragments of the template files we want to reuse in multiple places – with and without adaptions to the particular use case.

Template Functions			
▷ Definition 5.4.22. stpl python supplies the template functions			
1. include(((tpl)), ((vars))), where ((tpl)) is another template file and ((vars)) a set of variable declarations (for ((tpl))).			
2. defined($\langle\!\langle var \rangle\!\rangle$) for checking definedness $\langle\!\langle var \rangle\!\rangle$			
3. get($\langle\!\langle var \rangle\!\rangle$, $\langle\!\langle default \rangle\!\rangle$): return the value of $\langle\!\langle var \rangle\!\rangle$, or $\langle\!\langle default \rangle\!\rangle$.			
4. setdefault(《name》,《val》)			
▷ Example 5.4.23 (Including Header and Footer in a template). In a coherent web site, the web pages often share common header and footer parts. Realize this via the following page template:			
% include('header.tpl', title='Page Title') Page Content % include('footer.tpl')			
\triangleright Example 5.4.24 (Dealing with Variables and Defaults).			
% setdefault('text', 'No Text') <h1>{{get('title', 'No Title')}}</h1> {{ text }} % if defined('author'): By {{ author }} % end			
FAU : 174 2025-06-05 CONTRACTOR			

5.4.3 Completing the Contact Form

We are now equipped to finish the contact form example

We now come back to our worked HTML example: the contact form from above. Here is the current state:

Back to our Contact Form (Current State)

ī

\triangleright A contact form and messa	ge receipt	(communicate via HTTPs request)
contact4.html <title>Contact</title> <form action="contact-aft</td><td>ter.html"></form>	contact—after.html <title></title>	
<pre><in2>Please enter a mess <input <h3="" name="msg" ty=""/>Your e-mail addres <input name="addr" text"="" typ="" value="xx @ xx.cd </pre></td><td>age:</12>
/pe="/> s: e="text" le"/></in2></pre>	<form action="contact4.html"> <h2> Your message has been submitted!</h2></form>	
 <input <br="" type="submit"/> value="Send mess	sage"/>	 <input <br="" type="submit"/> value="Continue"/>
	0 /	
GET contact—after.html? msg=Hi;addr=foo@bar. $\leftarrow \rightarrow C \hat{\omega}$	de	GET contact.html
🥑 Getting Started 🗎 FAL	(←) → ℃ @	(i) file:///Users/l
Please enter a message:	🥑 Getting Started 🗎	FAU Services News Math
Your e-mail address: xx @ xx.de	Your messag	ge has been submitted!
Send message	Continue	
▷ Problem: The answer is a	a static HTML o	document independent of form data.
▷ Solution: Generate the a	nswer programm	natically using the form data. (up next)
FAU	175	2025-06-05

There are two great flaws in the current state of the contact form:

- 1. The "receipt page" contact—after.html is static and does not take the data it receives from the contact form into account. It would be polite to give some record on what happened. We can fix this using bottle WSGI using the methods we just learned.
- 2. Nothing actually happens with the message. It should be either entered into an internal message queue in a database0 or ticketing system, or fed into an e-mail to a sales person. As we do not have access to the first, we will just use a Python library to send an e-mail programmatically.



112

contact.py	contact-after.tpl
from bottle import route, run, debug,	Message submitted!
template, request, get	
<pre>@get('/contact—after.html')</pre>	Return Address:
def new_item():	${addr} $
$data = {$ 'msg': request.GET.msg.strip(),	
'addr': request.GET.addr.strip()}	
send—contact—email(addr,msg)	Message Sent:
return template('contact—after',**data)	{{msg}}
run(host="localhost", port=8080)	
FAU : 176	2025-06-05

Fortunately, the only remaining part: actually sending off an e-mail to the specified mailbox is very easy: using the smtplib library we just create an e-mail message object, and then specify all the components.



Once we have the e-mail message object msg, we open a "SMTP connection" s send the message via its send_message method and close the connection by s.quit()). Again, the Python library hides all the gory details of the SMTP protocol.

CHAPTER 5. WEB APPLICATIONS

Chapter 6

Front-end Technologies

We introduce three important concepts for building modern web front-ends for web applications:

- 1. Client-side computation: manipulating the browser DOM via JavaScript.
- 2. Cascading Stylesheets (CSS) for styling the layout of HTML (and XML).
- 3. The jQuery library: a symbiosis of JavaScript and CSS ideas to make JavaScript coding easier and more efficient.

6.1 Dynamic HTML: Client-side Manipulation of HTML Documents

We now turn to client-side computation:

One of the main advantages of moving documents from their traditional ink-on-paper form into an electronic form is that we can interact with them more directly. But there are many more interactions than just browsing hyperlinks we can think of: adding margin notes, looking up definitions or translations of particular words, or copy-and-pasting mathematical formulae into a computer algebra system. All of them (and many more) can be made, if we make documents programmable. For that we need three ingredients:

- i) a machine-accessible representation of the document structure, and
- *ii)* a program interpreter in the web browser, and
- *iii)* a way to send programs to the browser together with the document.

We will sketch the WWW solution to this in the following.

To understand client-side computation, we first need to understand the way browsers render HTML pages.

Background: Rendering Pipeline in browsers

- ▷ **Observation:** The nested markup codes turn HTML documents into trees.
- Definition 6.1.1. The document object model (DOM) is a data structure for the HTML document tree together with a standardized set of access methods.
- ▷ **Rendering Pipeline:** Rendering a web page proceeds in three steps
 - 1. the browser receives a HTML document,
 - 2. parses it into an internal data structure, the DOM,



The most important concept to grasp here is the tight synchronization between the DOM and the screen. The DOM is first established by parsing (i.e. interpreting) the input, and is synchronized with the browser UI and document viewport. As the DOM is persistent and synchronized, any change in the DOM is directly mirrored in the browser viewpoint, as a consequence we only need to change the DOM to change its presentation in the browser. This exactly is the purpose of the client side scripting language, which we will go into next.

6.1.1 JavaScript in HTML



The example above already shows a JavaScript command: document.write, which replaces the content of the <body> element with its argument – this is only useful for testing and debugging purposes.

Current web applications include simple office software (word processors, online spreadsheets, and presentation tools), but can also include more advanced applications such as project management, computer-aided design, video editing and point-of-sale. These are only possible if we carefully balance the effects of server-side and client-side computation. The former is needed for computational resources and data persistence (data can be stored on the server) and the latter to keep personal information near the user and react to local context (e.g. screen size).

Here are three browser level functions that can be used for user interaction (and finer debugging as they do not change the DOM).

Browser-level JavaScript functions: 1			
▷ Example 6.1.5 (Logging to the browser console).			
	🕞 🗘 Inspektor 🗵 Konsole		
	Ausgabe filtern		
	GET https://hnu-webeng-a		
	<pre>> GET https://hnu-webeng-a Hallo Welt1</pre>		
FAU	180	2025-06-05	

The function console.log writes its argument into the console of the web browser.

It is primarily used for debugging the source code of a web page.

Example 6.1.6. If we want to know whether a function square has been executed we add calls to console.log like this:

```
function square (n) {
   console.log ("entered function square with argument " + n);
   return (n * n);
   console.log ("exited function square with result " + n * n);
  }
```

In the console we can check whether the content contains e.g. entered function square and moreover whether argument and value are as expected.

Browser-level JavaScript functions: 2		
ho Example 6.1.7 (Raising a	a Popup).	
alert("Dynamic HTML for	IWGS!")	
	Dynamic HTML for IWGS!	
	ОК	

	FAU	181	2025-06-05	CO Same a datas de Served
The	function alert creates a popup	that contains the argument.		
	Browser-level JavaScr	ipt functions: 3		
	ho Example 6.1.8 (Asking	for Confirmation).		
	var returnvalue = confirm("Dynamic HTML for IWGS!")			
		Dynamic HTML for IWGS!		
		Cancel OK		
	FAU	182	2025-06-05	

The function confirm creates a popup that contains the argument and a confirmation/cancel button pair and returns the corresponding Boolean value.

If the user clicks on the confirmation button, the returned value will be false and true for the cancel button.

Example 6.1.9. You can play with this in the following frizzle:

```
<html>
<head>
   <title>confirm</title>
   <script src="./client-js/jquery-3.6.4.min.js" type="application/javascript"></script>
   <style>
       .emph{
           color: blue;
       }
       .code{
           font-size: 110%;
       }
   </style>
</head>
<body>
   <h2>Live Demo of the JavaScript <span class="code emph">confirm</span> Function</h2>
   <textarea id="output" style="width:400px">
   </textarea>
   <textarea id="code" style="width:400px; height:400px">
   </textarea>
   Click <button onclick="openPopup()">here</button> to execute
     the <span class="code">confirm</span> function again!
   <'p'>
       Show <button onclick="showCode()">source code</button>
   <script type="application/javascript">
       function openPopup(){
          console.log("executed openPopup function");
var output="";
```

```
var returnValue=confirm ("Hello World!");
           if(returnValue==true){
               output="You clicked the OK button!" + "(return value: " + returnValue + ")";
           } else {
               output="You clicked the Cancel button!" + "(return value: " + returnValue
+ ")";
           }
           console.log(output);
           $("#output").html(output);
           $("p").show();
       }
       openPopup();
       function showCode(){
           console.log("executed showCode function");
           var func=openPopup.toString();
           //alert(func);
           $("#code").html(func);
       }
   </script>
</body>
</html>
```

JavaScript is a client side programming language, that means that the programs are delivered to the browser with the HTML documents and is executed in the browser. There are essentially three ways of embedding JavaScript into HTML documents:

Embedding JavaScript into HT	ML		
\triangleright In a <script> element in HTML, e.g.</td><td></td><td></td><td></td></tr><tr><td><script type="text/javascript"> function sayHello() { console.log('h</td><td>tello IWGS!');</td><td></td><td></td></tr><tr><td></script>			
▷ External JavaScript file via a <script< p=""></script<>	> element with src attrib	ute:	
<script src="</td" type="text/javascript"></script>			

A related – and equally important – question is, *when* the various embedded JavaScript fragments are executed. Here, the situation is more varied



▷ JavaScript in a script element	ent: during page load:	(not in a function)
<script type="text/javasc</td><td>ript"></script>		

The first key concept we need to understand here is that the browser essentially acts as an user interface: it presents the HTML pages to the user, waits for actions by the user – usually mouse clicks, drags, or gestures; we call them events – and reacts to them.

The second is that all events can be associated to an element node in the DOM: consider an HTML anchor node, as we have seen above, this corresponds to a rectangular area in the browser window. Conversely, for any point p in the browser window, there is a minimal DOM element e(p)that contains p recall that the DOM is a tree. So, if the user clicks while the mouse is at point p, then the browser triggers a click event in e(p), determines how e(p) handles a click event, and if e(p) does not, bubbles the click event up to the parent of e(p) in the DOM tree.

There are multiple ways a DOM element can handle an event: some elements have default event handlers, e.g. an HTML anchor $\langle a href="\langle URI \rangle" \rangle$ will handle a click event by issuing a HTTP GET request for (URI). Other HTML elements can carry event handler attributes whose JavaScript content is executed when the corresponding event is triggered on this element.

Actually there are more events than one might think at first, they include:

- 1. Mouse events; click when the mouse clicks on an element (touchscreen devices generate it on a tap); contextmenu: when the mouse right-clicks on an element; mouseover / mouseout: when the mouse cursor comes over / leaves an element; mousedown / mouseup: when the mouse button is pressed / released over an element; mousemove: when the mouse is moved.
- 2. Form element events; submit: when the visitor submits a <form>; focus: when the visitor focuses on an element, e.g. on an *<input>*.
- 3. keyboard events; keydown and keyup: when the visitor presses and then releases the button.
- 4. Document events; DOMContentLoaded:- when the HTML is loaded and processed, DOM is fully built, but external resources like pictures and stylesheets may be not yet loaded. load: the browser loaded all resources (images, styles etc); beforeunload / unload: when the user is leaving the page.
- 5. resource loading events; onload: successful load, onerror: an error occurred.

Let us now use all we have learned in an example to fortify our intuition about using JavaScript to change the DOM.

Example: Changing Web Pages Programmatically > Example 6.1.10 (Stupid but Fun).

6.2. CASCADING STYLESHEETS

<body></body>			
<h2>A Pyramid</h2>			
<div id="pyramid"></div>		Eine Pyramide	
<pre><script #";<br="" type="text/javas
var char = ">var triangle = ""; var str = ""; for(var i=0;i<=10;i++ str = str + char; triangle = triangle + } var elem = document.g elem.innerHTML=trian</pre></td><td>cript">){ - str + " " etElementById("pyramid"); gle;</td><td># ### #### ##### ####### ######## ######</td><td></td></tr><tr><td></script> </pre>			
FAU	185	2025-06-05	

The HTML document in Example 6.1.10 contains an empty $\langle div \rangle$ element whose id attribute has the value pyramid. The subsequent script element contains some code that builds a DOM nodeset of 10 text and $\langle br \rangle$ nodes in the triangle variable. Then it assigns the DOM node for the $\langle div \rangle$ to the variable elem and deposits the triangle node-set as children into it via the JavaScript innerHTML method.

We see the result on the right of Example 6.1.10. It is the same as if the #-strings and $\langle br/\rangle$ sequence had been written in the HTML which at least for pyramids of greater depth would have been quite tedious for the author.

6.2 Cascading Stylesheets

In this section we introduce a technology of digital documents which naturally belongs into chapter 4: the specification of presentation (layout, colors, and fonts) for marked-up documents.

6.2.1 Separating Content from Layout

As the WWW evolved from a hypertext system purely aimed at human readers to a Web of multimedia documents, where machines perform added-value services like searching or aggregating, it became more important that machines could understand critical aspects web pages. One way to facilitate this is to separate markup that specifies the content and functionality from markup that specifies human-oriented layout and presentation (together called "styling"). This is what "cascading style sheets" set out to do.

Another motivation for CSS is that we often want the styling of a web page to be customizable (e.g. for vision impaired readers).

CSS: Cascading Style Sheets

- ▷ Idea: Separate structure/function from appearance.
- ▷ **Definition 6.2.1.** Cascading Style Sheets (CSS) is a style language that allows authors and users to attach style (e.g., fonts, colors, and spacing) to HTML and XML documents.
- ▷ **Example 6.2.2.** Our text file from ??? with embedded CSS:

<html> <head> <style type="text/css"> body {background-color:#d0e4fe;} h1 {color:orange; text-align:center;} p {font-family:"Verdana"; font-size:20px;} </style> </head> <body> <h1>CSS example</h1></body></html>	 CSS example Hello IWGS!.
	2025-06-05

Now that we have seen the example, let us fix the basic terminology of CSS.



In Example 6.2.5 the selectors are just element names, they specify that the respective declaration blocks apply to all elements of this name.

We explore this new technology by way of an example. We rework the title box from the HTML example above – after all treating author/affiliation information as headers is not very semantic. Here we use div and span elements, which are generic block-level (i.e. paragraph-like) and inline containers, which can be styled via CSS classes. The class titlebox is represented by the CSS selector .titlebox.

A Styled HTML Title Box (Source) ▷ Example 6.2.6 (A style Title Box). The HTML source: <head> <title>A Styled HTML Title</title> k rel="stylesheet" type="text/css" href="style.css"/>

6.2. CASCADING STYLESHEETS



And here is the result in the browser:

A Styled	HTML Title Box (Result)	
	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	
	Anatomy of a HTML Web Page Michael Kohlhase FAU ERLANGEN-NÜRNBERG	
Fau	: 189 2025-06-05	ि इत्यावनात्मात्मव्यव्यव्य

6.2.2 Worked Example: The Contact Form

To fortify our intuition on CSS, we take up the "contact form" example from above and improve the layout in a step-by-step process concentrating on one aspect at a time.

CSS in Practice: The Contact Form Example (Continued)
⊳ Recap: The unstyled contact form – Dream vs. Reality

<title>Contact</title> <form action="contact-after.html"> <h2>Please enter a message:</h2> <input name="msg" type="text"/> <h3>Your e-mail address:</h3> <input <br="" name="addr" type="text"/>value="xx @ xx.de"/> <input <br="" type="submit"/>value="Send message"/> </form>	Please type in a Message Your e-mail address: KX
$\begin{array}{c} \hline \\ \hline $	= FAI age: s:
Add a CSS file with font information link rel="stylesheet" type="text/css" href="csscontact1.css" /> <input class="important" message"="" send="" type="subn
value="/>	nit" Contact
<pre>body {font-size: 62.5%; font-family: "Trebuchet MS", "Arial", "Helvetica", "Verdana", "sans-serif"} .important{font-style: italic;} input[type="submit"]{font-weight: bold;}</pre>	Getting Started FAU
ightarrow Add lots of color	(ooops, what about the size)

6.2. CASCADING STYLESHEETS



This worked example should be enough to cover most layout needs in practice. Note that in most use cases, these generally layout primitives will have to be combined in different and may be even new ways.

Actually, the last "improvement" may have gone a bit overboard; but we used it to show how absolute positioning of images (or actually any CSS boxes for that matter) works in practice.

6.2.3 A small but useful Fragment of CSS

CSS is a huge ecosystem of technologies, which is spread out over about 100 particular specifications – see [CSSa] for an overview.

We will now go over a small fragment of CSS that is already very useful for web applications in more detail and introduce it by example. For a more complete introduction, see e.g. [CSSc]. Recall that selectors are the part of CSS rules that determine what elements a rule affects. We now give the most important cases for our applications.

CSS Selectors
▷ Question: Which elements are affected by a CSS rule?
▷ Elements of a given name (optionally with given attributes)
$\triangleright \text{ Selectors: name } \widehat{=} \langle \langle elname \rangle \rangle, \text{ attributes } \widehat{=} [\langle \langle attname \rangle \rangle = \langle \langle attval \rangle \rangle]$
\triangleright Example 6.2.7. p[xml:lang='de'] applies to
▷ Any element with a given class attributes
⊳ Selector: .((classname))
ightarrow Example 6.2.8important applies to <((el)) class='important'> ((el))
▷ The element with a given id attribute
\triangleright Selector: $\#\langle\!\langle id \rangle\!\rangle$
ightarrow Example 6.2.9. #myRoot applies to <((el)) id='myRoot'> ((el))
▷ Note: Multiple selectors can be combined in a comma separated list.
> For a full list see https://www.w3schools.com/cssref/css_selectors.asp.
FAU : 191 2025-06-05 EXECUTION

We now come to one of the most important conceptual parts of CSS: the box model. Understanding it is essential for dealing with CSS based layouts.





As a summary of the above, we can visualize the CSS box model in a diagram:



We now come to a topic that is quite mind-boggling at first: The "cascading" aspect of CSS stylesheets. Technically, the story is quite simple, there are two independent mechanisms at work:

- inheritance: if an element is fully contained in another, the inner (usually) inherits all properties of the outer.
- rule priorization: if more than one selector applies to an element (e.g. one by element name and one by id attribute), then we have to determine what rule applies.

Technically, priorization takes care of them in an integrated fashion.



CHAPTER 6. FRONT-END TECHNOLOGIES



But do not despair with this technical specification, you do not have to remember it to be effective with CSS practically, because the rules just encode very natural "behavior". And if you need to understand what the browser – which implements these rules – really sees, use the integrated page inspector tool (see slide 199 for details).

We now look at an example to fortify our intuition.

Cascading of selectors	in CSS: Priorization	Example	
▷ Example 6.2.12. Can you	explain the colors in the we	b browsers below?	
<h1>Layout with CSS<div class="<br" id="important">I am <span blue"="" class="marke
</div></td><td>1>
="> edimportant">very importa</div></h1>	nt		
.markedimportant {backgro #important {background— .blue {background—color:bl #important {background—c	und—color:red !important} color:green} ue} color:yellow}	←→ C ŵ Getting Started ► FAU ► Services Layout with CSS Law (Regularized)	
FAU	195	2025-06-05	COME PROFILIS RESERVED

For instance, the words "very important" get a red background, as the class markedimportant is marked as important by the CSS keyword !important, which makes (cf. rule 1 above) the color red win agains the color yellow inherited from the parent <div> element (rule 7 above). Let us now look at CSS inheritance in a little more detail.

Cascading in CSS: Inheritance	
Definition 6.2.13. Child elements can inherits some from their parents. In a nutshell:	e properties (called inheritable)
text-related properties are inheritable; e.g. color, f list—style, and text—align	ont, letter—spacing, line—height,
box-related properties are not; e.g. background height, width, margin, padding, position, and tex	l, border, display, float, clear, t—align.
▷ Note: Inheritance is integrated into priorization.	(recall case 7. above)
\rhd Inheritance makes for consistent text properties and	smaller CSS stylesheets.
	2025-06-05

So far, we have looked at the mechanics of CSS from a very general perspective. We will now come to a set of CSS behaviors that are useful for specifying layouts of pages and texts. Recall that CSS is based on the box model, which understands HTML elements as boxes, and layouts as properties of boxes nested in boxes (as the corresponding HTML elements are).

If we can specify how inner boxes float inside outer boxes - via the CSS float rules, we can already do quite a lot, as the following examples show.

128



One of the important applications of the content/form separation made possible by CSS is to tailor web page layout to the screen size and resolution of the device it is viewed on. Of course, it would be possible to maintain multiple layouts for a web page one per screensize/resolution class, but a better way is to have one layout that changes according to the device context. This is what we will briefly look at now.



6.2.4 CSS Tools

In this subsection we introduce a technology of digital documents which naturally As CSS has grown to be very complex and moreover, the browser DOM of which CSS is part can even be modified after loading the HTML (see ???), we need tools to help us develop effective and maintainable CSS.





In CSS we can specify colors by various names, but the full range of possible colors can only specified by numeric (usually hexadecimal) numbers. For instance in ???, we specified the background color of the page as #d0e4fe;, which is a pain for the author. Fortunately, there are tools that can help.



6.3 jQuery: Write Less, Do More

While JavaScript is fully sufficient to manipulate the HTML DOM, it is quite verbose and tedious to write. To remedy this, the web developer community has developed libraries that extend the

JavaScript language by new functionalities that more concise programs and are often used Instead of pure JavaScript.

jQuery: Write Less, Do More
Definition 6.3.1. JQuery is a feature-rich JavaScript library that simplifies tasks like HTML document traversal and manipulation, event handling, animation, and Ajax.
⊳ Using:
Download from https://jquery.com/download/, save on your system (remember where)
ightarrow integrate into your HTML (usually in the <head>)</head>
<script src="client-js/jquery-3.2.1.min.js" type="text/javascript"></script>
or from the internet directly (only works if you are online)
<script src="https://ajax.googleapis.com/ajax/libs/jquery/3.2.1/jquery.min.js"></script>
FAU : 201 2025-06-05 CONTRACTOR

The key feature of jQuery is that it borrows the notion of "selectors" to describe HTML node sets from CSS actually, jQuery uses the CSS selectors directly and then uses JavaScript-like methods to act on them. In fact, the name jQuery comes from the fact that selectors "query" for nodes in the DOM.



We will now show a couple of jQuery methods for inserting material into HTML elements and discuss their behavior in examples

Inserting Material into the DOM

6.3. JQUERY: WRITE LESS, DO MORE

▷ Inserting before the first	child:		
\$('#content').prepend(funct	tion(){return 'in front';});		
▷ Inserting after the last cl	nild:		
\$('#content').append('	Hello');		
\$('#content').append(funct	ion(){ return 'in the back'; }));	
► Inserting before/after an	olomont.		
	ciement.		
\$('#price').before('Price:');			
\$('#price').after(' EUR')			
FAU	203	2025-06-05	COME ALGENIS RESERVED

Let us fortify our intuition about dynamic HTML by going into a more involved example. We use the toggle method from the jQuery layout layer to change visibility of a DOM element. This method adds and removes a style="display:none" attribute to an HTML element and thus toggles the visibility in the browser window.

Applications and useful tricks in Dynamic HTML
▷ Observation: jQuery is not limited to adding material to the DOM.
\triangleright Idea: Use jQuery to change CSS properties in the DOM as well.
▷ Example 6.3.3 (Visibility). Hide document parts by setting CSS style attributes to display:none.
<html> <head> <title>Toggling</title> <style type="text/css">#dropper { display: none; }</style> <script src="https://ajax.googleapis.com/ajax/libs/jquery/3.2.1/jquery.min.js"></script> <script language="JavaScript" type="text/javascript"></script></head></html>



▷ The HTML has a butt	ton with text "hover".		
▷ The jQuery code select method.	cts it via its id and catches	its hover event via th	e hover()
Description This takes two function	ons as arguments:		
 ▷ The first is called v it leaves. ▷ The first changes of 	when the mouse moves into changes the button color to	the button, the second rever	nd when ts this.
FAU	205	2025-06-05	Some disting deserved
Fun with Buttons (Th	tree easy Interaction	s)	

\triangleright Example 6.3.5 (A Button that Uncovers Text).	
<div id="readPoint"></div>	
<button class="read" style="display:block">Read More</button>	
<pre><div id="rlext" style="display:none; width:200px; clear:left"></div></pre>	
A read-more button is not only a call-to-action, but it also organizes	rested
the will use the button $\langle \mathbf{hr} \rangle$	esteu,
<script type="text/javascript"></td><td></td></tr><tr><td><pre>\$(".read").click(function() {\$("#rText").toggle("slow",function(){\$(".read").to</pre></td><td>$ggle()$;});});</td></tr><tr><td></script>	
\triangleright The HTML has two buttons (one of them visible) and a text.	
▷ The jQuery code selects both buttons via their read class.	
\triangleright A click event activates the .click() method taking an event handler fun	ction:
\triangleright This selects the text via its id attribute rTeX and	
\triangleright uses the toggle() method which changes the display between none an	d block.
\sim The first parameter of togeta() is a duration for the animation	
I he second is a completion function to be run after animation finitial	shes.
▷ Here complection function makes the respective other button visi	ble (read
more/less).	
E 206 2025-06-05	SOMERDHISTISSEWED



▷ The HTML has a but	ton with text "sound" and a	an onclick attribute.	
\triangleright That activates the pla	ySound function on a URL:		
▷ The playSound function	n is defined in the script eler	ment: it	
\triangleright logs the action an	d URL in the browser conso	ole,	
⊳ makes a new audi	o object a, which		
ho plays it via the pla	y() method.		
= //			

For reference, here is the full code of the examples in one file:

```
<html>
<head>
  <title>Buttons</title>
  <script src="https://code.jquery.com/jquery-3.4.1.min.js" type="text/javascript"></script>
  <style type="text/css">
   button {color: white; font-size: large; background-color: blue;
                  width: 110px; height: 40px; border-radius: 20px;}
    div[id$="Point"] {display: inline-block;}
  </style>
</head>
<body>
  <h1 id="top">Look how easy interaction is ... </h1>
  <div id="hoverPoint">
    <button id="hover">hover</button>
    <script type="text/javascript">
       \label{eq:started} $$ $ "#hover" hover" (function () {$(this).css("background-color", "red");}, function () {$(this).css("background-color", "blue");}); $
    </script>
  </div>
  <div id="readPoint">
    <button class="read" style="display:block">Read More</button>
    <button class="read" style="display:none">Read Less</button>
    <div id="rText" style="display:none; width:200px; clear:left">
      A read-more button is not only a call-to-action, but it also organizes
      the screen area management in a non-wasteful way. If and only if users are interested,
      they will use the button. {<}br/{>}
    </div>
    <script type="text/javascript">
      $(".read").click(function() {$("#rText").toggle("slow",function(){$(".read").toggle()});});
    </script>
  </div>
  <div id="soundPoint">
    <button id="sound" onclick="playSound('laugh.mp3')">Sound</button>
    <script type="text/javascript">
     function playSound(url) {
        console.log("Call playSound with " + url);
        const a = new Audio(url);
        a.play();
        }
    </script>
  </div>
</body>
</html>
```

It has a bit more general CSS and includes jQuery in the beginning.
Chapter 7

Practical Aspects of Web Applications

7.1 Web Applications: Recap



Recap: Web Application Front-end







7.2 Maintaining State in Web Sites

There is one problem however with web applications that is difficult to solve with the technologies so far. We want web applications to give the user a consistent user experience even though they are made up of multiple web pages. In a regular application we we only want to log in once and expect the application to remember e.g. our username and password over the course of the various interactions with the system. For web applications this poses a technical problem which we now discuss.



Definition 7.2.3. Third-part	y cookies are used by a	dvertising companies to tra-	ck
users across multiple sites.	(but you can turn	off, and even delete cookie	s)
FAU	210	2025-06-05	\$5165Y/80

Note that both solutions to the state problem are not ideal, for usernames and passwords the URL-based solution is particularly problematic, since HTTP transmits URLs in GET requests without encryption, and in our example passwords would be visible to anybody with a packet sniffer. Here cookies are little better, since they can be requested by any web site you visit.

7.3 Access Control and Management

Now that we have a basic web application running, we can start adding features. The most important one is access control to restrict who can access more critical functionalities of the web application, such as deleting or updating database record.

There are many technologies for access control, many use advanced features like browser cookies. Here we want to introduce the simplest one: HTTP basic authentication is built into the fabric of the world wide web: it is part of the HTTP protocol that drives it.

As HTTP basic authentication is unsafe (it sends usernames and passwords over the network only lightly encrypted), we also add a discussion on how to upgrade the web application to HTTPS.

The full source is available at https://gl.mathhub.info/courses/FAU/IWGS/blob/master/ source/databases/code/books-app-https.py. The respective template files are siblings.

Access Control and Management			
▷ Problem: Anyone can write, edit, and delete records from the books database.			
▷ Solution: Implement a password-based log in procedure and restrict write/ed- it/delete access to logged-in agents.			
▷ Let's fix some terminology before we continue			
Definition 7.3.1. Access control (AC) is the selective restriction of access to a resource or place, access management describes the corresponding process.			
> Access management usually comprises both authentication and authorization.			
Definition 7.3.2. Authorization refers to a set of rules that determine who is allowed to do what with a collection of resources.			
For our books application we need four things			
1. a browser interaction to query the user for username and password			
2. a way to transport them to the web application program			
3. a method for checking the username/password (authentication)			
4. a way the specify who can do what. (authorization)			
Realization : 1./2. via HTTP, 4. via bottle basic auth, implement 3. directly.			
EAU : 211 2025-06-05			

HTTP basic authentication is a simple mechanism in the HTTP protocol that standardizes the transmission of username/password information the "handshake" that leads to its acquisition.



The message sequence diagram in Definition 7.3.3 shows the basic handshake process that establishes authentication and the delivery of restricted resources to an authenticated user.

The diagram shows the details of the communication between client and server (symbolized by the two vertical lines). The top arrow is a normal HTTP GET request (without a Authorization field).

But – as the resource that is requested is access-restricted – the server does not just answer with a HTTP "200 OK" and the resource, instead the server response with a HTTP "401 Unauthorized" code, which contains a description of the reason for the restriction.

When the browser receives the 401 response, it asks the user for a username and password e.g. with a popup form like the one shown in Definition 7.3.3, possibly displaying the reason string – here "private". This information is then send to the server in a second GET request, this time with the username/password information in the Authorization request.

The server checks the user/password data and – depending on the result of the check – sends a HTTP response "200 OK" together with the resource or a "403 Forbidden" (without the resource).

▲ One thing that we have not discussed here is that most browsers store the username/password information and supply it to the server – often directly in any outgoing requests – which makes it hard to test authentication and unauthenticated behavior in web application development. A useful trick here is – if you are logged into http://example.org – to address a GET request to http://abc@example.org. Background: HTTP basic authentication allows you to set user/password information directly by prepending 《user》:《pass》 to the authority of the URI used in a HTTP request.

Of course, HTTP basic authentication is supported by the bottle WSGI framework.

Basic Auth in Bottle
► Idea: Support the server side of HTTP basic authentication in bottle web-apps.
► Implementation: New decorator @auth_basic(《function》) to mark a route as password-protected.
► Usage: Decorate every route we want to restrict access of with @auth_basic(《function》), where 《function》 is a function that takes two string arguments (user name and password) and returns a Boolean for the authorization decision.

What happens behind the scene here is clear from the authentication handshake explained in ???

Basic Auth in Bottle: Minimal Viable Example			
▷ Example 7.3.4. A web application with restricted route.			
from bottle import run, get, auth_basic			
def check(user, password):			
return user == "miko" and password == "test"			
Øget("/")			
@auth_basic(check)			
def protected():			
return "Authorized access granted!"			
run(host="localhost", port=8000)			
\triangleright Idea: Mix restricted and open routes in a partially restricted application.			
Extension: Use different check functions for different levels of restriction (us roles)	er		



This was easy enough. But one problem remains: in HTTP basic authentication, usernames and passwords are not confidential when they are transported over the network. The simplest way to ensure confidentiality is to layer encryption on top of HTTP, which is just what the HTTPS protocol does.

7.4 HTTPS: Secure/Encrypted HTTP

HTTPS: HTTP over TLS				
Definition 7.4.1. Hypertext Transfer Protocol Secure (HTTPS) is an extension of the Hypertext Transfer Protocol (HTTP) for confidential communication over a computer network. HTTPS achieves this by running HTTP over a TLS connection.				
▷ Consequences for Web Applications: We can use HTTP as usual, except				
 ▷ we gain communication privacy and server authentication, ▷ server and browser need to speak HTTPS, (most do) ▷ the server needs a public key certificate and a private key. 				
In bottle, we can just swap out the HTTP server to one that can do HTTPS: run(host='localhost',port='8888', server='gunicorn',keyfile='key.pem',certfile='cert.pem')				
install it first with pip install gunicorn.				
Problem: Where to get the certificate file cert.pem and private key key.pem?				
EAU : 215 2025-06-05 EXTENSION				

For publically deploying a HTTPS based web application we need real TLS certificates. Fortunately, there is a relatively simple way of obtaining them.

Getting a Real TLS Certificate via Let's-Encrypt				
▷ Intuition: HTTPS is the new "regular HTTP" on the web!				
▷ Definition 7.4.2. In a public key infrastructure, the TLS certificate is issued by a certificate authority, an organization chartered to verify identity and issue TLS certificates.				
Commercial certificate authorities sell trust. (for a lot of money) They certify e.g. that the https://bmw.com is under control of BMW AG.				
Idea: Finding out that you have control over a particular web site on the web can be automated, if you run a program on the server host.				
Definition 7.4.3. Let's Encrypt is a not for profit certificate authority that does this and issues free TLS certificates. (to encourage HTTPS adoption)				
▷ Concretely: on a linux server you need two steps				
1. install certbot (usually via your package manager)				

2. then sudo /usr/local/bin/certbot certonly ——standalone will generate certs.
Details at https://letsencrypt.org.
▷ Success: ≥ 1.000.000.000 TLS certificates, 200.000.000 sites since 2016

We have only covered the basic ideas behind certificate authorities and Let's Encrypt here, but this should enable you to figure out the rest from the Let's Encrypt web site.

Chapter 8

What did we learn in IWGS-1?

Outline of IWGS 1:

▷ Programming in Python:

(main tool in IWGS)

2025-06-05

- ▷ Systematics and culture of programming
- ▷ Program and control structures
- ▷ Basic data structures like numbers and wordsstring, character encodings, unicode, and regular expressions
- ▷ Electronic documents and document processing:
 - ⊳ text files
 - ▷ markup systems, HTML, and CSS
 - ▷ XML: Documents are trees.

 \triangleright Web technologies for interactive documents and web applications

- ▷ internet infrastructure: web browsers and server
- ▷ server-side computation: bottle routing and
- ▷ client-side interaction: dynamic HTML, JavaScript, HTML forms

▷ Web application project (fill in the blanks to obtain a working web app)

217

Outline of IWGS-II:

▷ Databases

FAU

- ▷ CRUD operations, querying, and python embedding
- $_{\vartriangleright}$ XML and JSON for file based data storage
- \rhd $\operatorname{BooksApp}:$ a Books Application with persistent storage
- ▷ Image processing
 - \triangleright Basics



Bibliography

- [All18] Jay Allen. New User Tutorial: Basic Shell Commands. 2018. URL: https://www. liquidweb.com/kb/new-user-tutorial-basic-shell-commands/ (visited on 10/22/2018).
- [BLFM05] Tim Berners-Lee, Roy T. Fielding, and Larry Masinter. Uniform Resource Identifier (URI): Generic Syntax. RFC 3986. Internet Engineering Task Force (IETF), 2005. URL: http://www.ietf.org/rfc/rfc3986.txt.
- [CSSa] All CSS Specifications. URL: https://www.w3.org/Style/CSS/specs.en.html (visited on 01/12/2020).
- [CSSb] CSS Specificity. URL: https://en.wikipedia.org/wiki/Cascading_Style_ Sheets#Specificity (visited on 12/03/2018).
- [CSSc] CSS Tutorial. URL: https://www.w3schools.com/css/default.asp (visited on 12/02/2018).
- [DH98] S. Deering and R. Hinden. Internet Protocol, Version 6 (IPv6) Specification. RFC 2460. Internet Engineering Task Force (IETF), 1998. URL: http://www.ietf.org/ rfc/rfc2460.txt.
- [Ecm] ECMAScript Language Specification. ECMA Standard. 5th Edition. Dec. 2009.
- [ET] xml.etree.ElementTree The ElementTree XML API. URL: https://docs.python. org/3/library/xml.etree.elementtree.html (visited on 04/15/2021).
- [Fie+99] R. Fielding et al. Hypertext Transfer Protocol HTTP/1.1. RFC 2616. Internet Engineering Task Force (IETF), 1999. URL: http://www.ietf.org/rfc/rfc2616.txt.
- [Hic+14] Ian Hickson et al. HTML5. A Vocabulary and Associated APIs for HTML and XHTML. W3C Recommentation. World Wide Web Consortium (W3C), Oct. 28, 2014. URL: http://www.w3.org/TR/html5/.
- [HL11] Martin Hilbert and Priscila López. "The World's Technological Capacity to Store, Communicate, and Compute Information". In: Science 331 (2011). DOI: 10.1126/ science.1200970. URL: http://www.sciencemag.org/content/331/6018/692. full.pdf.
- [HWC] The Hello World Collection. URL: http://helloworldcollection.de/ (visited on 11/23/2018).
- [JKI] Jonas Betzendahl. juptyter.kwarc.info Documentation. URL: https://kwarc. info/teaching/IWGS/jupyter-documentation.pdf (visited on 08/29/2020).
- [Kar] Folgert Karsdorp. Python Programming for the Humanities. URL: http://www.karsdorp.io/python-course/ (visited on 10/14/2018).
- [Koh06] Michael Kohlhase. OMDoc An open markup format for mathematical documents [Version 1.2]. LNAI 4180. Springer Verlag, Aug. 2006. URL: http://omdoc.org/ pubs/omdoc1.2.pdf.

- [Koh08] Michael Kohlhase. "Using LATEX as a Semantic Markup Format". In: Mathematics in Computer Science 2.2 (2008), pp. 279-304. URL: https://kwarc.info/kohlhase/ papers/mcs08-stex.pdf.
- [LP] Learn Python Free Interactive Python Tutorial. URL: https://www.learnpython. org/ (visited on 10/24/2018).
- [LXMLa] *lxml XML and HTML with Python*. URL: https://lxml.de (visited on 12/09/2019).
- [LXMLb] *lxml API*. URL: https://lxml.de/api/ (visited on 12/09/2019).
- [LXMLc] The lxml.etree Tutorial.URL: https://lxml.de/tutorial.html (visited on 12/09/2019).
- [Nor+18a] Emily Nordmann et al. Lecture capture: Practical recommendations for students and lecturers. 2018. URL: https://osf.io/huydx/download.
- [Nor+18b] Emily Nordmann et al. Vorlesungsaufzeichnungen nutzen: Eine Anleitung für Studierende. 2018. URL: https://osf.io/e6r7a/download.
- [P3D] Python 3 Documentation. URL: https://docs.python.org/3/ (visited on 09/02/2014).
- [PyRegex] Rodolfo Carvalho. *PyRegex Your Python Regular Expression's Best Buddy*. URL: http://www.pyregex.com/ (visited on 12/03/2018).
- [Pyt] re Regular expression operations. online manual at https://docs.python.org/2/library/re.html. URL: https://docs.python.org/2/library/re.html.
- [Rfc] DOD Standard Internet Protocol. RFC. 1980. URL: http://tools.ietf.org/rfc/ rfc760.txt.
- [RHJ98] Dave Raggett, Arnaud Le Hors, and Ian Jacobs. HTML 4.0 Specification. W3C Recommendation REC-html40. World Wide Web Consortium (W3C), Apr. 1998. URL: http://www.w3.org/TR/PR-xml.html.
- [She24] Esther Shein. 2024. URL: https://cacm.acm.org/news/the-impact-of-ai-oncomputer-science-education/.
- [sTeX] sTeX: A semantic Extension of TeX/LaTeX. URL: https://github.com/sLaTeX/ sTeX (visited on 05/11/2020).
- [Sth] A Beginner's Python Tutorial. http://www.sthurlow.com/python/. seen 2014-09-02. URL: http://www.sthurlow.com/python/.
- [STPL] Simple Template Engine. URL: https://bottlepy.org/docs/dev/stpl.html (visited on 12/08/2018).
- [Swe13] Al Sweigart. Invent with Python: Learn to program by making computer games. 2nd ed. online at http://inventwithpython.com. 2013. ISBN: 978-0-9821060-1-3. URL: http: //inventwithpython.com.
- [Xam] apache friends Xampp. http://www.apachefriends.org/en/xampp.html. URL: http://www.apachefriends.org/en/xampp.html.

Appendix A

Excursions

As this course is predominantly an overview over (some) CS tools useful in the humanities and social sciences and not about the theoretical underpinnings, we give the discussion about these as a "suggested readings" chapter here.

A.1 Internet Basics

We will show aspects of how the internet can cope with this enormous growth of numbers of computers, connections and services. The growth of the internet rests on three design decisions taken very early on. The internet

- 1. is a packet-switched network rather than a network, where computers communicate via dedicated physical communication lines.
- 2. is a network, where control and administration are decentralized as much as possible.
- 3. is an infrastructure that only concentrates on transporting packets/datagrams between computers. It does not provide special treatment to any packets, or try to control the content of the packets.

The first design decision is a purely technical one that allows the existing communication lines to be shared by multiple users, and thus save on hardware resources. The second decision allows the administrative aspects of the internet to scale up. Both of these are crucial for the scalability of the internet. The third decision (often called "net neutrality") is hotly debated. The defenders cite that net neutrality keeps the Internet an open market that fosters innovation, where as the attackers say that some uses of the network (illegal file sharing) disproportionately consume resources.



FAU : 219 2025-06-05 EXERCISE

These ideas are implemented in the Internet Protocol Suite, which we will present in the rest of the section. A main idea of this set of protocols is its layered design that allows to separate concerns and implement functionality separately.

The Intenet Protocol Suite			
۵	Definition A.1.2. The Inte Protocol Suite (commonly known TCP/IP) is the set of communicat protocols used for the internet other similar networks. It structu into 4 layers.	rnet Layer Application Layer Transport Layer Internet Layer Link Layer	e.g. HTTP, SSH UDP, TCP IPv4, IPsec Ethernet, DSL
۵	Layers in TCP/IP: TCP/IP uses sulation to provide abstraction of cols and services. An application (the highest level model) uses a set of protocols to s data down the layers, being furth capsulated at each level.	encap- proto- of the send its her en- Frame Frame data	Data Application UDP data Transport ta Internet Frame footer Link
FAU	: 220	2025-	© 06-05
The Internet as a Network of Networks			
	nario). Consider a situation with two internet host computers communicate across local network boundaries.	Network Conne Host \rightarrow Router \rightarrow Router	ections er $Host$ B
⊳	Network boundaries are consti- tuted by internetworking gateways (routers).	Stack Connec	
	Definition A.1.4. A router is a purposely customized com- puter used to forward data among computer networks beyond directly connected devices.	Transport	et Internet
	A router implements the link and internet layers only and has two network connections.	Ethernet Satellite, etc.	Ethernet
			œ

We will now take a closer look at each of the layers shown above, starting with the lowest one.

Instead of going into network topologies, protocols, and their implementation into physical signals that make up the link layer, we only discuss the devices that deal with them. Network Interface

A.1. INTERNET BASICS

controllers are specialized hardware that encapsulate all aspects of link-level communication, and we take them as black boxes for the purposes of this course.

Network Interfaces				
▷ The nodes in the internet are computers, the edges communication channels				
▷ Definition A.1.5. A network interface controller (NIC) is a hardware device that handles an interface to a computer network and thus allows a network-capable device to access that network.				
Definition A.1.6. Each NIC contains a unique number, the media access control address (MAC address), identifies the device uniquely on the network.				
MAC addresses are usually 48-bit numbers issued by the manufacturer, they are usually displayed to humans as six groups of two hexadecimal digits, separated by hy- phens (-) or colons (:), in transmission order, e.g. 01-23-45-67-89-AB, 01:23:45:67:89:AB				
	Definition A.1.7. A network interface	Layer	e.g.	
	is a software component in the operat-	Application Layer	HTTP, SSH	
\triangleright	ing system that implements the higher	Transport Layer	ТСР	
	levels of the network protocol (the NIC	Internet Layer	IPv4, IPsec	
	handles the lower ones).	Link Layer	Ethernet, DSL	
ightarrow A computer can have more than one network interface. (e.g. a router)				
Fau	: 222	2025-0	6-05	

The next layer ist he Internet Layer, it performs two parts: addressing and packing packets.



Internet Protocol and IP Addresses

e

2025-06-05

- Definition A.1.11. The internet mainly uses Internet Protocol Version 4 (IPv4) [Rfc], which uses 32 bit numbers (IPv4 addresses) for identification of network interfaces of computers.
- ightarrow IPv4 was standardized in 1980, it provides 4,294,967,296 (2³²) possible unique addresses. With the enormous growth of the internet, we are fast running out of IPv4 addresses.
- ▷ Definition A.1.12. Internet Protocol Version 6 [DH98] (IPv6), which uses 128 bit numbers (IPv6 addresses) for identification.
- Although IP addresses are stored as binary numbers, they are usually displayed in human-readable notations, such as 208.77.188.166 (for IPv4), and 2001:db8:0:1234:0:567:1:1 (for IPv6).

FAU : 224

The internet infrastructure is currently undergoing a dramatic retooling, because we are moving from IPv4 to IPv6 to counter the depletion of IP addresses. Note that this means that all routers and switches in the internet have to be upgraded. At first glance, it would seem that this problem could have been avoided if we had only anticipated the need for more the 4 million computers. But remember that TCP/IP was developed at a time, where the internet did not exist yet, and it's precursor had about 100 computers. Also note that the IP addresses are part of every packet, and thus reserving more space for them would have wasted bandwidth in a time when it was scarce.

We will now go into the detailed structure of the IP packets as an example of how a low-level protocol is structured. Basically, an IP packet has two parts: the "header", whose sequence of bytes is strictly standardized, and the "payload", a segment of bytes about which we only know the length, which is specified in the header.

The Structure of IP Packets

 \triangleright **Definition A.1.13.** IP packets are composed of a 160b header and a payload. The IPv4 packet header consists of:

b	name	comment
4	version	IPv4 or IPv6 packet
4	Header Length	in multiples 4 bytes (e.g., 5 means 20 bytes)
8	QoS	Quality of Service, i.e. priority
16	length	of the packet in bytes
16	fragid	to help reconstruct the packet from fragments,
3	fragmented	$DF \stackrel{\sim}{=} "Don't fragment"/MF \stackrel{\sim}{=} "More Fragments"$
13	fragment offset	to identify fragment position within packet
8	TTL	Time to live (router hops until discarded)
8	protocol	TCP, UDP, ICMP, etc.
16	Header Checksum	used in error detection,
32	Source IP	
32	target IP	
	optional flags	according to header length

 \triangleright Note that delivery of IP packets is not guaranteed by the IP protocol.

FAU

2025-06-05

A.1. INTERNET BASICS

As the internet protocol only supports addressing, routing, and packaging of packets, we need another layer to get services like the transporting of files between specific computers. Note that the IP protocol does not guarantee that packets arrive in the right order or indeed arrive at all, so the transport layer protocols have to take the necessary measures, like packet re-sending or handshakes,



We will see that there are quite a lot of services at the network application level. And indeed, many web-connected computers run a significant subset of them at any given time, which could lead to problems of determining which packets should be handled by which service. The answer to this problem is a system of "ports" (think pigeon holes) that support finer-grained addressing to the various services.



On top of the transport-layer services, we can define even more specific services. From the perspective of the internet protocol suite this layer is unregulated, and application-specific. From a user perspective, many useful services are just "applications" and live at the application layer.

The Application Layer				
 Definition A.1.18. The application layer of the internet protocol suite contains all protocols and methods that fall into the realm of process-to-process communications via an Internet Protocol (IP) network using the Transport Layer protocols to establish underlying host-to-host connections. Example A.1.19 (Some Application Layer Protocols and Services). 				
	BitTorrent	Peer-to-peer	Atom	Syndication
	DHCP	Dynamic Host Configuration	DNS	Domain Name System
	FTP	File Transfer Protocol	HTTP	HyperText Transfer
	IMAP Internet Message Access IRCP Internet F		Internet Relay Chat	
	NFS	Network File System	NNTP	Network News Transfer
	NTP	Network Time Protocol	POP	Post Office Protocol
	RPC	Remote Procedure Call	SMB	Server Message Block
	SMTP	Simple Mail Transfer	SSH	Secure Shell
	TELNET	Terminal Emulation	WebDAV	Write-enabled Web
EAU : 228 2025-06-05 EAU				

The domain name system is a sort of telephone book of the internet that allows us to use symbolic names for hosts like kwarc.info instead of the IP number 212.201.49.189.



Let us have a look at a selection of the top-level domains in use today.

Domain Name Top-Level Domains

- \rhd .com ("commercial") is a generic top-level domain. It was one of the original top-level domains, and has grown to be the largest in use.
- \triangleright .org ("organization") is a generic top-level domain, and is mostly associated with non-profit organizations. It is also used in the charitable field, and used by the open-

source movement. Government sites and Political parties in the US have domain names ending in .org

- ▷ .net ("network") is a generic top-level domain and is one of the original top-level domains. Initially intended to be used only for network providers (such as Internet service providers). It is still popular with network operators, it is often treated as a second .com. It is currently the third most popular top-level domain.
- ▷ .edu ("education") is the generic top-level domain for educational institutions, primarily those in the United States. One of the first top-level domains, .edu was originally intended for educational institutions anywhere in the world. Only postsecondary institutions that are accredited by an agency on the U.S. Department of Education's list of nationally recognized accrediting agencies are eligible to apply for a .edu domain.

FAU

Domain Name Top-Level Domains

- \triangleright .info ("information") is a generic top-level domain intended for informative website's, although its use is not restricted. It is an unrestricted domain, meaning that anyone can obtain a second-level domain under .info. The .info was one of many extension(s) that was meant to take the pressure off the overcrowded .com domain.
- .gov ("government") a generic top-level domain used by government entities in the United States. Other countries typically use a second-level domain for this purpose, e.g., .gov.uk for the United Kingdom. Since the United States controls the .gov Top Level Domain, it would be impossible for another country to create a domain ending in .gov.
- ▷ .biz ("business") the name is a phonetic spelling of the first syllable of "business". A generic top-level domain to be used by businesses. It was created due to the demand for good domain names available in the .com top-level domain, and to provide an alternative to businesses whose preferred .com domain name which had already been registered by another.
- ▷ .xxx ("porn") the name is a play on the verdict "X-rated" for movies. A generic top-level domain to be used for sexually explicit material. It was created in 2011 in the hope to move sexually explicit material from the "normal web". But there is no mandate for porn to be restricted to the .xxx domain, this would be difficult due to problems of definition, different jurisdictions, and free speech issues.

FAU

231

2025-06-05

C

2025-06-05

Note: Anybody can register a domain name from a registrar against a small yearly fee. Domain names are given out on a first-come-first-serve basis by the domain name registrars, which usually also offer services like domain name parking, DNS management, URL forwarding, etc.

The telnet Protocol
▷ Problem: We need a way to remotely operate networked computers via a shell.
▷ Idea: Send shell instructions and responses as text messages between a terminal

2025-06-05

client (a program on the local host) and a terminal server (a program on the remote host).

- ▷ Definition A.1.24. The telnet protocol uses TCP directly to send text based messages two networked computers. It customarily uses port 25.
- ⊳ **Remark**:

telnet is one of the oldest protocols in the TCP/IP protocol suite. It is no longer used much by itself (it is superseded by rsh and ssh), but still serves as a basis for other protocols, e.g. HTTP.

FAU

232

The next application-level service is the SMTP protocol used for sending e-mail. It is based on the telnet protocol for remote terminal emulation which we do not discuss here.

A Protocol Example: SMTP over telnet				
Definition A.1.25. The Simple Mail Transfer Protocol (SMTP) is a communica- tion protocol for electronic mail transmission based on telnet.				
▷ Example A.1.26. The SMTP protocol starts out by establishing identity				
▷ We call up the telnet service on the Jacobs mail server telnet exchange.jacobs-university.de 25				
\triangleright it identifies itself	(have some patience, it is very busy)			
Trying 10.70.0.128 Connected to exchange.jacobs-university.de. Escape character is '^]'. 220 SHUBCASO1.jacobs.jacobs-university.de Microsoft ESMTP MAIL Service ready at Tue. 3 May 2011 13:51:23 +0200				
▷ We introduce ourselves politely	(but we lie about our identity)			
helo mailhost.domain.tld				
▷ It is really very polite.				
250 SHUBCAS04.jacobs.jacobs-university.de Hello [10.222.1.5]				
	3 2025-06-05			



158

A.1. INTERNET BASICS

FAU



Essentially, the SMTP protocol mimics a conversation of polite computers that exchange messages by reading them out loud to each other (including the addressing information). We could go on for quite a while with understanding one Internet protocol after each other, but this is beyond the scope of this course (indeed there are specific courses that do just that). Here we only answer the question where these protocols come from, and where we can find out more about them.

235

Internet Standardization
▷ Question: Where do all the protocols come from?(someone has to manage that)
Definition A.1.29. The Internet Engineering Task Force (IETF) is an open stan- dards organization that develops and standardizes internet standards, in particular the TCP/IP and Internet protocol suite.
ightarrow All participants in the IETF are volunteers (usually paid by their employers)
Rough Consensus and Running Code: Standards are determined by the "rough consensus method" (consensus preferred, but not all members need agree) IETF is interested in practical, working systems that can be quickly implemented.
▷ Idea: running code leads to rough consensus or vice versa.
Definition A.1.30. The standards documents of the IETF are called Request for Comments (RFC). (more than 6300 so far; see http://www.rfceditor.org/)
EAU : 236 2025-06-05 CONTRACTOR

2025-06-05