General Information & Communication Technology 350101 GenICT I & II 2015 Partial Lecture Notes

Michael Kohlhase

School of Engineering & Science Jacobs University, Bremen Germany m.kohlhase@jacobs-university.de

September 20, 2015

Preface

This Document

This document contains the course notes for the those parts of the course General Information & Communication Technology I & II held at Jacobs University Bremen in the academic year 2014.

Contents: The document mixes the slides presented in class with comments of the instructor to give students a more complete background reference.

Caveat: This document is made available for the students of this course only. It is still a draft and will develop over the course of the current course and in coming academic years.

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.

Course Concept

Aims: The course 350101 "General Information & Communication Technology I/II" (GenICT) is a two-semester course that introduces concepts of Computer Science Concepts to non-CS students. The course is co-taught by four Jacobs Computer Science Faculty each covering a quarter of the materials.

Course Contents

Goal: We want to demonstrate both theoretical foundations of CS as Computer Science, and we want to provide practical knowledge helping students to cope with understanding and handling Computers, electronic documents and data, and the Web. Roughly the first half of the first semester is devoted to theoretical foundations and core concepts (Kohlhase and Jaeger), and the second half of the semester to the practical real-world stuff (Schnwlder and Baumann). Throughout the semester, students will be introduced stepwise to one of the main programming languages of today, Python.

Acknowledgments

Materials: The presentation of the programming language python uses materials prepared by Dr. Heinrich Stamerjohanns and Dr. Florian Rabe for the ESM Phython modules.

GenICT Students: The following students have submitted corrections and suggestions to this and earlier versions of the notes: Kim Philipp Jablonski, Tom Wiesing.

ii

Contents

	Preface	ii
	This Document	ii
	Course Concept	ii
	Course Contents	ii
	Acknowledgments	ii
1	Outline of the Course	1
2	Administrativa	3
т	ConICT 1 Modulo 1: Programming & Documents	5
I	Genici 1, module 1. i rogramming & Documents	0
1 3	Introducction to Programming	9
1 3	Introducction to Programming	9 9
1 3	Introducction to Programming 3.1 Introduction to Programming 3.2 Programming in Python	9 9 11
1 3	Introducction to Programming 3.1 Introduction to Programming 3.2 Programming in Python 3.3 Basic Data Structures	9 9 11 17
1 3	Introducction to Programming 3.1 Introduction to Programming 3.2 Programming in Python 3.3 Basic Data Structures 3.3.1 Numbers	9 9 11 17 18
1 3	Introducction to Programming 3.1 Introduction to Programming 3.2 Programming in Python 3.3 Basic Data Structures 3.3.1 Numbers 3.3.2 Characters & Strings	9 9 11 17 18 21
1 3 4	Introducction to Programming 3.1 Introduction to Programming 3.2 Programming in Python 3.3 Basic Data Structures 3.3.1 Numbers 3.3.2 Characters & Strings Computing with Documents	9 9 11 17 18 21 27
1 3 4 5	Genice 1 1, Module 1. Trogramming & Documents Introducction to Programming 3.1 Introduction to Programming 3.2 Programming in Python 3.3 Basic Data Structures 3.3.1 Numbers 3.3.2 Characters & Strings Computing with Documents Structured and Web Documents	9 9 11 17 18 21 27 31
1 3 4 5	Introducction to Programming 3.1 Introduction to Programming 3.2 Programming in Python 3.3 Basic Data Structures 3.3.1 Numbers 3.3.2 Characters & Strings Computing with Documents Structured and Web Documents 5.1 Multimedia Documents on the World Wide Web	9 9 11 17 18 21 27 31 31

Chapter 1 Outline of the Course

This course gives a broad (and in a few places, also a bit more in-depth) introduction to Computer Science for non-CS students. We want to demonstrate both theoretical foundations of CS as Computer Science, and we want to provide practical knowledge helping students to cope with understanding and handling Computers, electronic documents and data, and the Web. Roughly the first half of the semester is devoted to theoretical foundations and core concepts, and the second half of the semester to the practical real-world stuff. Throughout the semester, students will be introduced stepwise to one of the main programming languages of today, Python.

Outline of the C	ourse I: Theory			
⊳ Introduction to Inf	formation Processing	(Michael	Kohlhase)	
▷ Introduction to	Programming in general and	python as the course	e language	
⊳ Encoding and o	computing numbers on the co	mputer		
▷ Encoding and t "texts")	ranslating characters and doc	cuments (also	known as	
⊳ Describing the XML)	structure of programs and do	cuments (regular ex	(pressions,	
▷ Legal foundation	ons of Intellectual Property	(Copyr	ight, etc.)	
⊳ From Symbols to	Computing	(Herbe	rt Jaeger)	
 Boolean logic (what "bits" are and how to combine them to to "digital computation") 				
▷ Complexity: so others are fast	me computer programs take v – must this be so?	very long to load or	complete,	
CUMBER CHINESES WAT	©: Michael Kohlhase	1		





Chapter 2

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 about Computer Science concepts as efficient and painless as possible.



Even though the lecture itself will be the main source of information in the course, there are various resources from which to study the material.

Textbooks, H	landouts and Inform	mation, Forum		
⊳ No required t	extbook, but course notes,	posted slides		
▷ Course notes de/teaching,	in PDF will be posted at ht /GenICTFall2015	ttp://minds.jacobs-un	iversity.	
▷ Everything wi	II be posted on PantaRhei	(Notes+assignments+cou	urse forum)	
⊳ announcer	nents, contact information,	course schedule and caler	ndar	
b discussion among your fellow students(careful, I will occasionally check for academic integrity!)				
ho http://pa	anta.kwarc.info	(use your Ja	cobs login)	
Set Up Pantal	Rhei Access: to get notifica	ations		
⊳ 1) Log into	PantaRhei,			
2) find the o	course GenICT Fall 2015	,		
3) request n	nembership	(I will ap	pprove you)	
\triangleright if there are pr	oblems send e-mail me (m	kohlhase@jacobs-unive	rsity.de)	
SOMME FILENTIS FIESTERVED	©: Michael Kohlhase	5	JACOBS UNIVERSITY	

No Textbook: There is no single textbook that covers the course. Instead we have a comprehensive set of course notes (this document). They are provided in two forms: as a large PDF that is posted at the course web page and on the PantaRhei system. The latter is actually the preferred method of interaction with the course materials, since it allows to discuss the material in place, to play with notations, to give feedback, etc. The PDF file is for printing and as a fallback, if the PantaRhei system, which is still under development, develops problems.

But of course, there is a wealth of literature on the subject, and the references at the end of the lecture notes can serve as a starting point for further reading. We will try to point out the relevant literature throughout the notes.

Part I

GenICT 1, Module 1: Programming & Documents

We start off the course by giving a very brief introduction to programming. We use python as the main programming language. As many students already have experience through the python lab, we will only briefly recap – and introduce to those students who have no experience yet – the python language

Chapter 3

Introducction to Programming

To understand programming we need to understand a bit about computers – the devices programs run on first

3.1 Introduction to Programming

Programming is an important and distinctive part of "Information and Communication Technology" – 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 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".



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

CHAPTER 3. INTRODUCCTION TO PROGRAMMING



hardware-software-programming

10



Program Execution

6

▷ Algorithm: informal description of what to do (good enough for humans)





▷ Program: formal version of the algorithm (needed for computers)

- ▷ Interpreter: reads a program and executes it directly
 - ▷ special case: interactive interpretation (lets you experiment easily)
- ▷ Compiler: translates a program (the source) into another program (the binary) in a much simpler language for optimized execution on hardware directly.
- Deservation 3.1.2 Compilers are more efficient, but more cumbersome for development

©: Michael Kohlhase

8

JACOBS UNIVERSIT

3.2 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 GenICT.

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



Installing python: python can be installed from http://python.org \sim "Downloads", as a Windows installer or a MacOSX disk image. For linux it is best installed via the package manager.

The download will install the python interpreter and the python shell IDLE3 that can be used for interacting with the interpreter directly. The latter also comes with an integrated editor for writing python programs. This editor gives you python syntax highlighting, and interpreter and debugger integration. In short, IDLE3 is an integrated development environment for python.

Arithmetic Expressions in python





Comments in python

 \triangleright It is highly advisable to insert comments into your programs

 \triangleright Single line comments start with a #

3.2. PROGRAMMING IN PYTHON

▷ Multiline comments start and end with three quotes (single or double: """ or · · ·) > Comments are ignored by python but are useful information for the programmer JACOBS

12

©: Michael Kohlhase

Variables in python

```
\triangleright Idea: values (of expressions) can be given a name for later reference
▷ Definition 3.2.1 A variable is a storage location which contains a value and
  an associated identifier – the variable name.
\triangleright A variable name can be used in expressions everywhere its value could be.
\triangleright in python:
    \triangleright = declares variable name and assigns value.
    ▷ variable names start with letter or _, cannot be keywords (case-sensitive)
              >>> foot = 30.5
              >>> inch = 2.54
              >>> 6 * foot + 2 * inch
              188.08
              >>> 3 * Inch
              Traceback (most recent call last):
                File "<pyshell#3>", line 1, in <module>
                   3 * Inch
              NameError: name 'Inch' is not defined
              >>>
▷ Example 3.2.2 (Swapping Variables)
  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)
                                                                            V JACOBS
UNIVERSITY
e
                       (C): Michael Kohlhase
                                                          13
```

Data Types in python

- \triangleright programs process data, which can be combined by operators
- \triangleright Data types group data into types
 - ▷ 1, 2, 3, etc. are data of type "integer"

▷ hello is data of type "string"

 \triangleright Data types determine which operators can be applied

 \triangleright In python, every value has a type, variables can have any type

▷ Definition 3.2.3 python has the following three basic data types

Data type	Name	Example
Integers	int	1, -5, 0
Floats	float	1.2, .125, -1.0
Strings	str	"Hello", 'Hello', "123", 'a'

 \triangleright The type of a variable is automatically determined when first used

```
firstVariable = 23 # integer
weight = 3.45 # float
first = 'Hello' # str
```

▷ Observation 3.2.4 python is strongly typed, i.e. types have to match

▷ Use datatye conversion functions int(), float(), and str() to adjust types

```
> Example 3.2.5 >>> 3+"hello"
Traceback (most recent call last):
File "<pyshell#1>", line 1, in <module>
3+"hello"
TypeError: unsupported operand type(s) for +: 'int' and 'str'
>>> str(4)+"hello"
'4Hello'
```

©: Michael Kohlhase

14

V JACOBS

Functions in python

▷ Observation: sometimes programming tasks are repetitive

```
print("Hello Peter, how are you today? How about some GenICT?")
print("Hello Roxana, how are you today? How about some GenICT?")
print("Hello Frodo, how are you today? How about some GenICT?)
...
```

 \triangleright Idea: We can automate the repetitive part by functions

ho Definition 3.2.6 A python function is defined by a code snippet of the form

```
def f(p_1, \ldots, p_n):

"""{docstring, what does this function do on parameters

:param p_i: document arguments}

"""

{this is the code in the function}

{more code, it can contain p_1, \ldots, p_n, and f}

return {value (e.g. text or number) of the function call}

{some code outside the function}
```

3.2. PROGRAMMING IN PYTHON





Branching in python

▷ Definition 3.2.12 Branching via if/then/else statements





Lists in python

- ▷ Definition 3.2.14 A list is a finite enumeration of objects
- \rhd In python lists can be written as a list of comma-separated expressions between square brackets.
- ▷ Example 3.2.15 (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 3.2.16 List elements can be assessed by specifying ranges
>>> list1[0]

```
'physics'
>>> list1[-2]
1997
>>> list2[1:4]
[2, 3, 4]

▷ Example 3.2.17 lists can be constructed by python functions
>>> list(range(1,6,2))
[1,3,5]
range(1,6,2) makes a range object from 1 to 6 with step 2, list makes a
list from it.
```



Now that we have acquired some basic programming skills, we will drill into how the basic data structures – numbers and strings – are actually represented in the computer. This will give us some insights on how to deal with them in practice.

3.3 Basic Data Structures

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

3.3.1 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.



18

3.3. BASIC DATA STRUCTURES

|--|

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 mankind, but it had tremendous consequences. There is archaeological evidence that in ancient times, namely at least some 8000 to 10000 years ago, men 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.



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 encoding.

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 encoding a meaning by mapping them into the unary natural numbers via a mapping φ_b . This is the the same process we did for the logics; we are now doing it for number

19

systems. However, here, we also want to ensure that the meaning mapping φ_b is a bijection, since we want to define the arithmetics on the encodings by reference to The arithmetical operators on the unary natural numbers.

Commonly Used Positional Number Systems							
⊳ Exam	$ ho { m Example ~ 3.3.7}$ The following positional number systems are in common use.						
	name	set	base	digits	example]	
	unary	\mathbb{N}_1	1	/	////1]	
	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	$162098_{10} \text{ or } 162098$		
	hexadecimal	\mathbb{N}_{16}	16	0,1,,9,A,,F	$FF3A12_{16}$		
decimal) Trick: Group triples or quadruples of binary digits into recognizable chunks(add leading zeros as needed)							
$\triangleright \triangleright 110001101011100_2 = \underbrace{0110_2}_{6_{16}} \underbrace{0011_2}_{6_{16}} \underbrace{0101_2}_{5_{16}} \underbrace{1100_2}_{C_{16}} = 635C_{16}$							
$\triangleright 110001101011100_2 = \underbrace{110_2 \ 001_2 \ 101_2 \ 011_2 \ 100_2}_{101_2 \ 011_2 \ 100_2} = 61534_8$							
$\succ F3A_{16} = \underbrace{F_{16}}_{1111_2} \underbrace{3_{16}}_{0011_2} \underbrace{A_{16}}_{1010_2} = 111100111010_2, 4721_8 = \underbrace{4_8}_{100_2} \underbrace{7_8}_{111_2} \underbrace{2_8}_{010_2} \underbrace{1_8}_{011_2} = 100111010001_2$							
CC Some rights reserved		©: M	chael K	ohlhase	24		

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 we (compound) digits for them. 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 circuits, since binary digits can be represented physically by current/no current.

Arithmetics in Positional Numb	per Systems
▷ For arithmetics just follow elementery s ▷ Tom Lehrer's "New Math"	chool rules (for the right base)
⊳ Example 3.3.9	
Addition base 4	binary multiplication 1 0 1 0
$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$
	$1 \ 1 \ 1 \ 1 \ 0 \ 0$

3.3. BASIC DATA STRUCTURES

C: Michael Kohlhase	25	
---------------------	----	--

3.3.2 Characters & Strings

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 subsection 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 a being able to encode a large set of characters and the representational capabilities in the time of punch cards (see below).



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.

A Punchcard

- ▷ A punch card is a piece of stiff paper that contains digital information represented by the presence or absence of holes in predefined positions.
- \triangleright Example 3.3.11 This punch card encoded the FORTRAN statement Z(1) = Y + W(1)

STATEMENT 2	1	FORTRAN STA	TEMENT	S S SFICATION
		0 0		0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
I 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	2222222222222222	2222222222222222222222222	222222222222222222222222222222222222222	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3	333333333333333	3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3	3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3	3 3 3 3 3 3 3 3 3 3 3 3 3 3
4444444444	444444444444444444444444444444444444444	444444444444444444444444444444444444444	4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4	4 4 4 4 4 4 4 4 4 4 4 4 4
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6	555555555555555555555555555555555555555	1 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	55555555555555 6666666 8 6666
רדרדקרדרקר,	111111111111111	111111111111111111111111111111111111111	111111111111111111111111111111111111111	11111
8 8 8 8 8 8 8 8 8 8	8 8 8 8 8 8 8 8 8 8 8 8			88888888888888
99999999999999 11234567891011 III234567891011	9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9	9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9	9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9	9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

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.



▲ Note: Example 3.3.12 and Example 3.3.13 (or any other examples in this lecture) is not

production code, but didactially motivated – to show you what you can do with the objects we are presenting in python.

In parcticular, if we "lowercase" a character that is already lowercase – e.g. by lc('c'), then we get out of the range of the ASCII code: the answer is x83, which is the character with the hexadecimal code 83 (decimal 130).

In production code (e.g. the python lower method), we would have some range checks, etc.

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 a byte (eight bits) is the preferred unit in information technology. Therefore there have been a whole zoo of extensions, 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)					
▷ Problem: Many European ch	aracters are not repre	esented (e.g. è,ñ	,ü,β,)		
⊳ European ASCII Variants: Ex	\triangleright European ASCII Variants: Exchange less-used characters for national ones				
$ \begin{tabular}{lllllllllllllllllllllllllllllllllll$					
▷ Definition 3.3.15 (ISO-Latin (ISO/IEC 8859)) 16 Extensions of ASCII to 8-bit (256 characters) ISO-Latin 1 \doteq "Western European", ISO-Latin 6 \doteq "Arabic", ISO-Latin 7 \doteq "Greek"					
▷ Problem: No cursive Arabic, Asian, African, Old Icelandic Runes, Math,					
Idea: Do something totally different to include all the world's scripts: For a scalable architecture, separate					
▷ what characters are available from the (character set)					
▷ bit string-to-character mapping (character encoding)					
©: Micha	el Kohlhase	29			

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.





Note that there is indeed an issue with space-efficient encoding 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 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					
Definition 3.3.20 A character UCS code points.	ter encodir	ig is a map	ping from	bit strings	to	
Idea: Unicode supports multiple encodings (but not character sets) for effi- ciency						
\triangleright Definition 3.3.21 (Unicod	de Transfo	ormation	Format)			
⊳ UTF-8, 8-bit, variable-widt ASCII.	th encoding	, which ma>	kimizes con	npatibility v	vith	
⊳ UTF-16, 16-bit, variable-w	vidth encodi	ng	(p	opular in A	sia)	
⊳ UTF-32, a 32-bit, fixed-wi	dth encodin	g		(for safe	ety)	
ho Definition 3.3.22 The UTF	-8 encoding	follows the	following e	encoding sch	neme	
Unicode	Byte1	Byte2	Byte3	Byte4	7	
U + 000000 - U + 00007F	0xxxxxxx					
U+000080 - U+0007FF	110xxxxx	10xxxxxx				

 $1110 \times \times \times$

11110xxx

10xxxxxx

10xxxxxx

10xxxxxx

10xxxxxx

10xxxxxx

U + 000800 - U + 00FFFF

U+010000 - U+10FFFF

3.3. BASIC DATA STRUCTURES

\triangleright Example 3	.3.23 $= U+0024$ is encoded as 0	0100100	(1 byte)
$\mathbf{c} = \mathbf{U} + 00A2$	(two bytes)		
$\in = U + 20A0$	(three bytes)		
SOME RIGHTS RESERVED	©: Michael Kohlhase	31	

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

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

Programming with UniCode strings is particularly simple, strings in python are UTF-8-encoded UniCode strings and all operations on them are UniCode-based² This makes the introduction to UniCode in python very short, we only have to know how to produce non-ASCII characters – which are on regular keyboards.

Unicode in python			
\triangleright the python str d	ata type is UniCode encod	led as UTF-8.	
⊳ How to write Uni	Code characters?: there a	re four ways	
▷ write them in your editor (make sure that it uses UTF-8)			
▷ otherwise use python escape sequences (try it!)			(try it!)
>>> "GREE '\u0394' >>> "\u0394"	K CAPITAL LETTER DEL # Using a 16-bit he	<pre>FA}" # Using the cl x value</pre>	haracter name
'\u0394'			
>>> "\U00000394" # Using a 32-bit hex value '\u0394'			
SOME FIGHTS RESERVED	©: Michael Kohlhase	32	

 $^{^{2}}$ Older programming languages have ASCII strings only, and UniCode strings are supplied by external modules.

Chapter 4 Computing with Documents

In this chapter we introduce methods to automatically deal with documents – actually large strings for the moment. We introduce "regular expressions", a domain-specific language for locating substrings of a particular form in a document. Regular expressions are useful in many documentrelated tasks, e.g. advanced searching and replacing, therefore most programming languages – python is no exception – integrate them as a sublanguage.

There are several dialects of regular expression languages that differ in details, but share the general setup and syntax. Here we introduce the UNIX variant.

Regular Expressions, see [RE]				
Definition 4.0.1 A regular expression (also called regexp) is a formal expression that specifies a set of strings.				
> Definition 4.0.2 (Meta-Characters for Regexps)				
	char	denotes		
		any single character		
	^	beginning of a string		
	\$	end of a string		
	[]	any single character in the brackets		
	[^]	any single character not in the brackets		
	()	marks a group		
	$\setminus n$	the n^{th} group		
		disjunction		
	*	matches the preceding element zero or more times		
	+	matches the preceding element one or more times		
	?	matches the preceding element zero or one times		
	$\{n,m\}$	matches the preceding element between n and m times		
	∖s	whitespace character		
	∖S	non-whitespace character		
All other characters match themselves, to match e.g. a ?, escape with a $\:\$				
SOME RIGHTS RESERVED		©: Michael Kohlhase 33	JACOBS UNIVERSITY	

Let us now fortify our intuition with some (simple) examples.

Re	Regular Expression Examples				
\triangleright	▷ Example 4.0.3 (Regular Expressions and their Values)				
	regexp	values			
	car	car			
	.at	cat, hat, mat,			
	[hc]at	cat, hat			
	[^c]at	hat, mat,(but not cat)			
	^[hc]at	hat, cat, but only at the beginning	of the line		
	[0-9]	Digits			
	[1-9][0-9]*	natural numbers			
	(.*)\1	mama, papa, wakawaka			
	cat dog	cat, dog			
 A regular expression can be interpreted by a regular expression processor (a program that identifies parts that match the provided specification) or a compiled by a parser generator. Example 4.0.4 (A more complex example) The following regexp times in a variety of formats, such as 10:22am, 21:10, 08h55, and 7.15 pm. 					
^(?:([0]?\d 1[012]) (?:1[3-9] 2[0-3]))[.:h]?[0-5] \d(?:\s?(?(1)(am AM pm PM)))?\$					
SOMERIC	ED INTERSECTION	©: Michael Kohlhase	34		

As we have seen regular expressions can become quite cryptic and long (cf. e.g. Example 4.0.4), so we need help in developing them. One way is to use one of the many regexp testers online

<pre>Iaying with Regular Expressions > If you want to play with regexps, go e.g. to http://regex101.com @@@@</pre>				
G regex101.com/#pyth		TATENO T	<u>م</u> (۲	
regular expressions 101	>_ regex tester 🛆 community 🥥 irc	INIBASE * Plats *	regex101 \$ donate	🚽 contact 🛕 bug reports & suggestions
FLWOR prec (php) payascript payascript TOOLS format regox (req c/) code generator # regex debugger pay past to community VESION COMTROL Pay ave regex ACCOUNT +) log in SETTINKS	REGULAR EXPRESSION "[Etb]at TEST STRING the rat bit the Cat	10000	EXPLANATION	character present in the list below in the list (1b literally (case sensitive) ters at literally (case sensitive) characted. Ittern matched anything in the
 coupley initiaspace wrap long lines colorize syntax use dark theme use minimal view 			QUICK REFERENCE PULL REFERENCE Q ★ most used tokens all tokens CATEGORIES ⊙ general tokens	MOST USED TOKENS A single character of: a, b or c [abc] A character except: a, b or c [vabc] A character in the range: a z [a:2] A character not in the range:
	SUBSTITUTION	0	🙏 anchors	A character in the range: a [a-zA-Z]

©:Michael Kohlhase 35	©: Michael Kohlhase 35	JACOBS UNIVERSITY
-----------------------	------------------------	----------------------

The **sed** stream editor is an example of a standalone utility – it is shipped with most operating systems – that uses regular expressions. It can be used to automate repetitive editing operations on files.



Example 4.0.8 shows the power of sed in combination with other utilities. Here we use the UNIX find utility that searches a file system for files with certain characteristics – here file names that match the regexp .*\.html.*txt— and executes the sed script cleanse we defined earlier.





We will now see what we can do with regular expressions in a practical example.

```
Example: Correcting and Anonymizing Documents
 \triangleright Example 4.0.9 We write an a program that makes simple corrections on
   documents and also crosses out all names.
    ▷ The worst president of the US, arguably was George W. Bush, right?
    ▷ However, are you famILIar with Paul Erdős or Henri Poincaré?(Unicode)
   Here is the program:
    ▷ we first initialize and load modules
      #!/usr/bin/env python3
      import re
      import sys
    \triangleright then we decode the argument and put it into a variable
      s = sys.argv[1]
    \triangleright We put put a space after a comma, (use r string prefix for "raw strings")
      s = re.sub(r", (\S)", r", \1", s)
    ▷ capitalize the first letter of a new sentence,
      s = re.sub(r"([\.\?!]) w*(\S)",
                   lambda (m):m.group(1),r" ".upper()+m.group(2), s)
    \triangleright next we make abbreviations for regular expressions to save space
      c = "[A-Z]"
      l = "[a-z]"
    ▷ remove capital letters in the middle of words
      s = re.sub("({1})({c}+)({1})"
          .format(l=l, c=c),
          lambda (m):"{0}{1}{2}".format(m.group(1), m.group(2).lower(),
          m.group(3)), s)
```

 \triangleright and we cross-out for official public versions of government documents,



Chapter 5

Structured and Web Documents

5.1 Multimedia Documents on the World Wide Web

We have seen the client-server infrastructure of the WWWeb, which essentially specifies how hypertext documents are retrieved. Now we look into the documents themselves.

In ?character-encodings? have already discussed how texts can be encoded in files. But for the rich documents we see on the WWWeb, we have to realize that documents are more than just sequences of characters. This is traditionally captured in the notion of document markup.

Document Markup					
Definition 5.1.1 (Document Markup) Document markupmarkup is the process of adding codes (special, standardized character sequences) to a docu- ment to control the structure, formatting, or the relationship among its parts.					
\triangleright Example	$ ho {f Example 5.1.2}$ A text with markup codes (for printing)				
	nus page, veries or radio]ANGELICA []]Angelica archangelcia Janes M. Stephens [] Jimes Reman z 26]INTRODUCTION [14 Jimes Reman Angelica is a Europeen perennial pl grown in this country as a culinary member of the pareley family, relat grows in fields and damp places fro Delaware and west to Minnesota.	12 Jimae Roman L.[24 Jimae Roman Ant sometimes herb. This ied to carrots, om Labrador to			
SOME AUGUST RESERVED	©: Michael Kohlhase	39			

There are many systems for document markup ranging from informal ones as in Definition 5.1.1 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.

WWWeb documents have a specialized markup language 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.



HTML was created in 1990 and standardized in version 4 in 1997. Since then there has HTML has been basically stable, even though the WWWeb 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. Acknowledging the growing discrepancy, the W3C has started the standardization of version 5 of HTML.



As the WWWeb evolved from a hypertext system purely aimed at human readers to an 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				
ightarrow Idea: Separate structure/function from ap	pearance.			
Definition 5.1.7 The Cascading Style Sheets (CSS), is a style sheet language that allows authors and users to attach style (e.g., fonts and spacing) to structured documents. Current version 2.1 is defined in [BCHL09].				
$ ho \ \mathbf{Example} \ 5.1.8$ Our text file from Examp	ple 5.1.5 with embedded CSS			
<html> <head>a <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>CCSS example</h1></body></html>	<pre> file:html > + - = = file:html > + - = = file:///Users/kohihase/tmp/test.html = C >> = Services = News = MathWeb = Jacobs = >> CSS example Hello GenCSII!. </pre>			
Hello GenCSII!. 				
,				
C: Michael Kohlhase	42 V Inversity			

5.2 Web Applications

In this section we show how with a few additions to the basic WWWeb infrastructure introduced in ?www-basics?, we can turn web pages into web-based applications that can be used without having to install additional software.

The first thing we need is a means to send information back to the web server, which can be used as input for the web application. Fortunately, this is already forseen by the HTML format.

HTML Forms: Submitting Information to the Web Server
$ ho {f Example 5.2.1}$ Forms contain input fields and explanations.
<form action="html_form_submit.asp" method="get" name="input"> Username: <input name="user" type="text"/> <input type="submit" value="Submit"/> </form>
The result is a form with three elements: a text, an input field, and a submit button, that will trigger a HTTP GET request to the URL specified in the action attribute.
Username: Submit

Image: State of the s	
---	--

As the WWWeb 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		
ho Idea: Why write HTML pages if we can also program them! (easy to do)		
Definition 5.2.2 A server-side scripting framework is a web server extension that generates web pages upon HTTP GET requests.		
Example 5.2.3 per1 is a scripting language with good string manipulation facilities. per1 CGI is an early server-side scripting framework based on this.		
Server-side scripting frameworks allow to make use of external resources (e.g. databases or data feeds) and computational services during web page generation.		
Problem: Most web page content is static (page head, text blocks, etc.) (and no HTML editing support in program editors)		
Idea: Embed program snippets into HTML pages. (only execute these, copy rest)		
Definition 5.2.4 A server-side scripting language is a server side scripting framework where web pages are generated from HTML documents with em- bedded program fragments that are executed in context during web page gen- eration.		
ho Note: No program code is left in the resulting web page after generation (important security concern)		
©: Michael Kohlhase 44		

To get a concrete intuition on the possibilities of server-side scripting frameworks, we will present PHP, a commonly used open source scripting framework. There are many other examples, but they mainly differ on syntax and advanced features.

PHP, a Server-Side Scripting Language
Definition 5.2.5 PHP (originally "Programmable Home Page Tools", later "PHP: Hypertext Processor") is a server-side scripting language with a C-like syntax. PHP code is embedded into HTML via special "tags" php and ?
$ ho \operatorname{\mathbf{Example}}$ 5.2.6 The following PHP program uses echo for string output
<html></html>
<pre><body><?php echo 'Hello world';?></body> </pre>
·,
$ ho {f Example 5.2.7}$ We can access the server clock in PHP (and manipulate it)
php</td

5.2. WEB APPLICATIONS

```
$tomorrow = mktime(0,0,0,date("m"),date("d")+1,date("Y"));
 echo "Tomorrow is ".date("d. m. Y", $tomorrow);
  ?>
 This fragment inserts tomorrow's date into a web page
\triangleright Example 5.2.8 We can generate pages from a database (here MySQL)
 <?php
 $con = mysql_connect("localhost","peter","abc123");
 if (!$con)
  {
   die('Could not connect: ' . mysql_error());
   }
 mysql_select_db("my_db", $con);
 $result = mysql_query("SELECT * FROM Persons");
 while($row = mysql_fetch_array($result))
   {
   echo $row['FirstName'] . " " . $row['LastName'];
   echo "<br />";
   }
 mysql_close($con);
\triangleright Example 5.2.9 We can even send e-mail via this e-mail form.
 <html><body>
  <?php
 if (isset($_REQUEST['email']))//if "email" is filled out, send email
   {//send email
   $email = $_REQUEST['email'] ;
   $subject = $_REQUEST['subject'] ;
   $message = $_REQUEST['message'] ;
   mail("someone@example.com", $subject,
   $message, "From:" . $email);
   echo "Thank you for using our mail form";}
 else //if "email" is not filled out, display the form
   {echo "<form method='post' action='mailform.php'>
    Email: <input name='email' type='text' /><br />
    Subject: <input name='subject' type='text' /><br />
    Message:<br />
    <textarea name='message' rows='15' cols='40'>
    </textarea><br />
    <input type='submit' />
    </form>";}
 2>
 </body></html>
                                                                            JACOBS
45
                      (C): Michael Kohlhase
```

What would we do in python

▷ Example 5.2.10 (HTML Hello World in python)

```
print("</html>")
print("<body>Hello world</body>")
print("</html>")
```

 \triangleright Why is this not as good as PHP?

▷ If HTML markup dominate, want to use a HTML editor (mode)



With server-side scripting frameworks like PHP, we can already build web applications, which we now define.

Web Applications: Usi	ng Applications w	vithout Installing		
Definition 5.2.11 A web application is a website that serves as a user inter- face for a server-based application using a web browser as the client.				
ho Example 5.2.12 Commo	nly used web application	s include		
 ▷ http://ebay.com; aud ▷ http://www.weather. ▷ http://slashdot.org ▷ http://github.com; s 	tion pages are generated com; weather information ; aggregation of news fea ource code hosting and p	from databases n generated weather fee eds/discussions project management	eds	
Common Traits: pages generated from databases and external feeds, content submission via HTML forms, file upload				
Definition 5.2.13 A web application framework is a software framework for creating web applications.				
$ ho \ Example \ 5.2.14$ The LAMP stack is a web application framework based on linux, apache, MySQL, and PHP.				
$ hinspace {f Example 5.2.15}$ A variant of the LAMP stack is available for Windows as XAMPP [XAM].				
©: Mic	hael Kohlhase	47	IACOBS UNIVERSITY	

Indeed, the first web applications were essentially built in this way. Note however, that as we remarked above, no PHP code remains in the generated web pages, which thus "look like" static web pages to the client, even though they were generated dynamically on the server.

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 login 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.





Note that 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 as cookies, since they can be requested by any website you visit.

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 documents. We will sketch the WWWeb solution to this in the following.

Dynamic HTML

- > Observation: The nested, markup codes turn HTML documents into trees.
- Definition 5.2.19 The document object model (DOM) is a data structure for the HTML document tree together with a standardized set of access methods.
- ▷ Note: All browsers implement the DOM and parse HTML documents into it; only then is the DOM rendered for the user.
- \triangleright Idea: generate parts of the web page dynamically by manipulating the DOM.
- Definition 5.2.20 JavaScript is an object-oriented scripting language mostly used to enable programmatic access to the DOM in a web browser.



Let us fortify our intuition about dynamic HTML by going into a more involved example.

Applications and useful tricks in Dynamic HTML
$ ho \ \mathbf{Example} \ 5.2.22$ hide document parts by setting CSS <code>style</code> attribs to <code>display:nor</code>
<html> <head> <style type="text/css">#dropper { display: none; }</style> <script language="JavaScript" type="text/javascript"> window.onload = function toggleDiv(element){ if(document.getElementById(element).style.display == 'none') {document.getElementById(element).style.display = 'block'} else if(document.getElementById(element).style.display == 'block') {document.getElementById(element).style.display == 'block') {document.getElementById(element).style.display = 'none'}} </script> </head> more Now you see it!</html>
Application: write "gmail" or "google docs" as JavaScript enhanced web applications. (client-side computation for immediate reaction) Image: Current Megatrend: Computation in the "cloud", browsers (or "apps") as user interfaces
©: Michael Kohlhase 50

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).

40

Bibliography

- [BCHL09] Bert Bos, Tantek Celik, Ian Hickson, and Høakon Wium Lie. Cascading style sheets level 2 revision 1 (CSS 2.1) specification. W3C Candidate Recommendation, World Wide Web Consortium (W3C), 2009.
- [ECM09] ECMAScript language specification, December 2009. 5th Edition.
- [RE] re regular expression operations. online manual at https://docs.python.org/2/ library/re.html.
- [RHJ98] Dave Raggett, Arnaud Le Hors, and Ian Jacobs. HTML 4.0 Specification. W3C Recommendation REC-html40, World Wide Web Consortium (W3C), April 1998.
- [XAM] apache friends xampp. http://www.apachefriends.org/en/xampp.html.