UNIB20005
# Language and Computation
Text Processing

Steven Bird, Department of Computing and Information Systems

# Accessing Text

- Sources on the web and on disk

- NLTK corpora: special collections, already prepared

- What if you want to study some other text?

- Today: issues with processing raw text

- Wednesday: how to do it in Python

# Text on the Web: Linguistic Value

| Google hits | *adore* | *love* | *like* | *prefer* |
|---|---|---|---|---|
| *absolutely* | 289,000 | 905,000 | 16,200 | 644 |
| *definitely* | 1,460 | 51,000 | 158,000 | 62,600 |
| ratio | 198:1 | 18:1 | 1:10 | 1:97 |

- Google queries, e.g.: "absolutely prefer"

- Benefits: coverage, ease of use

- Shortcomings of relying on a search engine: search patterns, inconsistencies, reproducibility, duplication

- Therefore: obtain texts and work with them directly (i.e. make your own corpus)

# Text on the Web: additional material

# Text on the Web: HTML source

```
     !doctype html public "-//W3C//DTD HTML 4.0 Transitional//EN" "http://www.w3.org/TR/REC-html40/loose.dtd">
     <html>
<head>
<title>BBC NEWS | Health | Blondes 'to die out in 200 years'</title>

...650 lines...

<font face="sans-serif" size="1"><span class="date">Friday, 27 September, 2002, 11:51 GMT 12:51 UK
</span></font>
     <div class="headlinestory"><b>Blondes 'to die out in 200 years'</b><br></div>
     <div class="inlineimage">
          <img height="180" vspace="0" border="0" width="300" alt=" " src="/media/images/38280000/jpg/_38280456_blonde300.jpg">
               <div class="caption"><font size="1">Scientists believe the last blondes will be in Finland</font><br></div>
     </div>
     <font class="body" face="sans-serif" size="2">
     <div class="bodytext">
     The last natural blondes will die out within 200 years, scientists believe.
<P>
A study by experts in Germany suggests people with blonde hair are an endangered species and will become extinct by 2202.
<P>
Researchers predict the last truly natural blonde will be born in Finland - the country with the highest proportion of blondes.
<P>
<P>
<!-- GENInlineBOX -->
     <table bgcolor="#FFFFCC" class="boxbody" cellspacing="0" width="150" border="0" cellpadding="3" align="right">
<!-- GENInlineQUOTE -->
     <tr><td><img src="/nol/shared/img/startquote.gif" width="23" height="18" border="0" valign="TOP" alt=""><br><div
class="boxbody">
     The frequency of blondes may drop but they won't disappear
     </div><img align="RIGHT" src="/nol/shared/img/endquote.gif" width="23" height="18" border="0" valign="ABSBOTTOM" alt=""><br
clear="ALL"></td></tr>
<!-- GENInlineNAME -->
     <tr><td bgcolor="cccc99"><div class="boxhead">
     Prof Jonathan Rees, University of Edinburgh
     </div></td></tr>
     </table>
But they say too few people now carry the gene for blondes to last beyond the next two centuries.
<P>
The problem is that blonde hair is caused by a recessive gene.
<P>
```

# Text on the Web: Extracting Text from HTML
*Simple method: delete all markup and collapse whitespace*

BBC NEWS | Health | Blondes 'to die out in 200 years' NEWS   SPORT   WEATHER   WORLD SERVICE   WHERE I LIVE -->   A-Z INDEX    SEARCH   You are in: Health   News Front Page Africa Americas Asia-Pacific Europe Middle East South Asia UK Business Entertainment Science/Nature Technology Health Medical notes ------------- Talking Point ------------- Country Profiles In Depth ------------- Programmes ------------- SERVICES Daily E-mail News Ticker Mobile/PDAs ------------- Text Only Feedback Help EDITIONS Change to UK Friday, 27 September, 2002, 11:51 GMT 12:51 UK Blondes 'to die out in 200 years' Scientists believe the last blondes will be in Finland **The last natural blondes will die out within 200 years, scientists believe. A study by experts in Germany suggests people with blonde hair are an endangered species and will become extinct by 2202. Researchers predict the last truly natural blonde will be born in Finland - the country with the highest proportion of blondes.** The frequency of blondes may drop but they won't disappear Prof Jonathan Rees, University of Edinburgh **But they say too few people now carry the gene for blondes to last beyond the next two centuries. The problem is that blonde hair is caused by a recessive gene. In order for a child to have blonde hair, it must have the gene on both sides of the family in the grandparents' generation. Dyed rivals The researchers also believe that so-called bottle blondes may be to blame for the demise of their natural rivals. They suggest that dyed-blondes are more attractive to men who choose them as partners over true blondes.** Bottle-blondes like Ann Widdecombe may be to blame **But Jonathan Rees, professor of dermatology at the University of Edinburgh said it was unlikely blondes would die out completely "Genes don't die out unless there is a disadvantage of having that gene or by chance. They don't disappear," he told BBC News Online. "The only reason blondes would disappear is if having the gene was a disadvantage and I do not think that is the case. "The frequency of blondes may drop but they won't disappear."** See also: 28 Mar 01 | Education What is it about blondes? 09 Apr 99 | Health Platinum blondes are labelled as dumb 17 Apr 02 | Health Hair dye cancer alert Internet links: University of Edinburgh The BBC is not responsible for the content of external internet sites Top Health

# Extracting text from HTML

- This is a non-trivial task!

- We will use a built-in "library function" to do this

# Unicode: Working at the level of Characters

# Unicode: Code Points and Glyphs

# Unicode (cont)



- Code pages: http://www.unicode.org/charts/
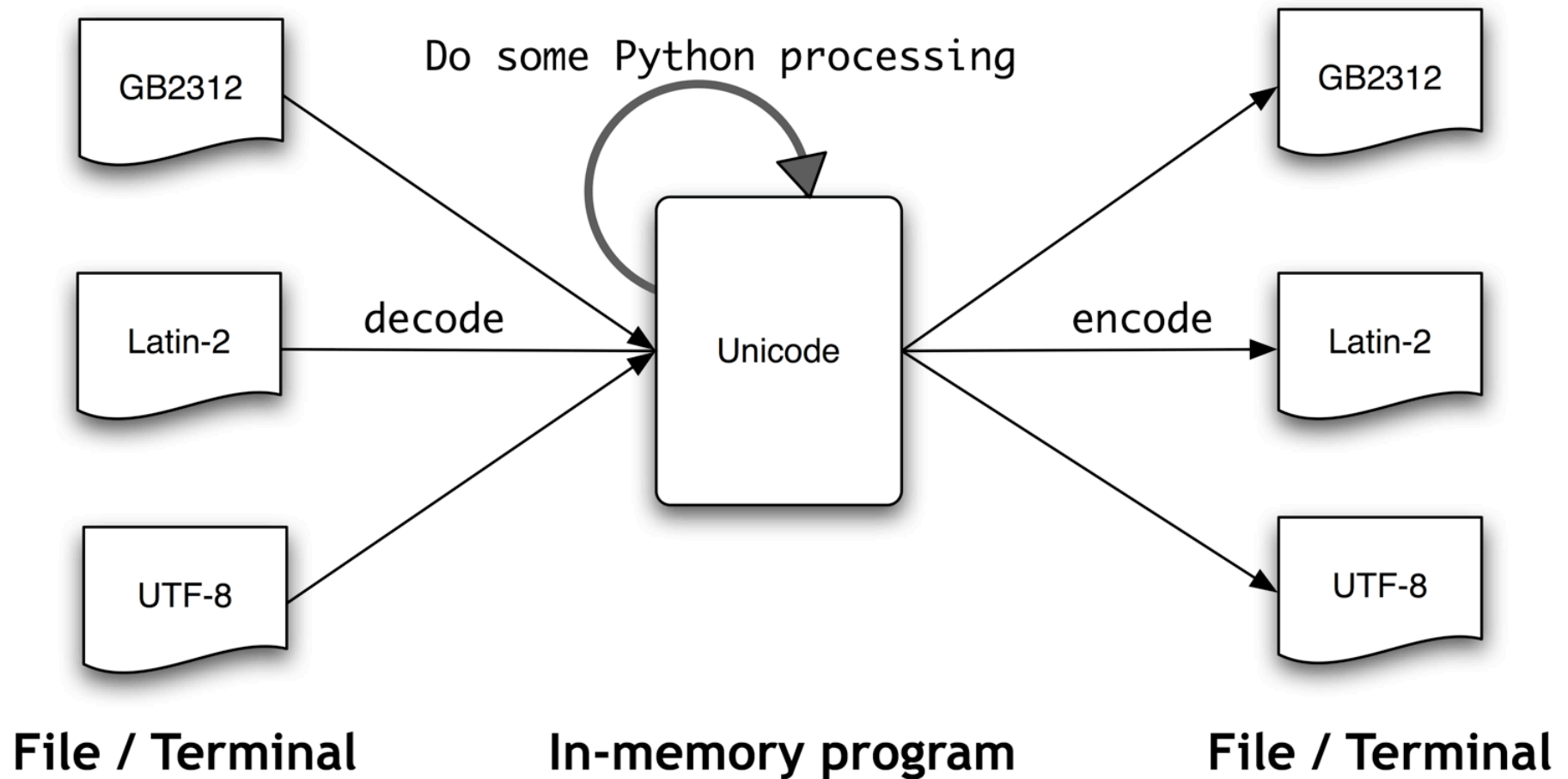
# E.g. Bengali

# Unicode Normalization

- typographic ligatures, e.g. ﬃ

- precomposed diacritics, e.g. o + ¨ = ö

- U+006F U+0308 = U+00F6

# UTF-8 Encoding

| 1st Byte | 2nd Byte | 3rd Byte | Number of Free Bits | Maximum Expressible Unicode Value |
|----------|----------|----------|---------------------|-----------------------------------|
| 0xxxxxxx |          |          | 7                   | 007F hex (127)                    |
| 110xxxxx | 10xxxxxx |          | (5+6)=11            | 07FF hex (2047)                   |
| 1110xxxx | 10xxxxxx | 10xxxxxx | (4+6+6)=16          | FFFF hex (65535)                  |

- UTF = "Unicode Transformation Format", e.g. UTF-8, UTF-16, UTF-32

- An *encoding* is how we represent a *codepoint* as a unique sequence of *bytes*

- For codepoints 0..127 we use one byte

- For codepoints 128..2047 we use two bytes

- For codepoints 2048..65535 we use three bytes

- This is called a "variable length" encoding

# Text Processing with Unicode

# Extracting characters from bytes

- This is a non-trivial task!

- We will use a built-in "library function" to do this

# Tokenization: From Strings to Tokens

# Strings

- Basic data type: sequence of characters

- what we get when we read from a file or URL

- we cannot process a text till we split it into tokens

# Tokenization

- Simple approach: split on whitespace:
  
  The last natural blondes will die out within 200 years, scientists believe.

- Split off punctuation as well:
  
  "The frequency of blondes may drop but they won't disappear."

- Harder case (*Alice in Wonderland*)
  
  'When I'M a Duchess,' she said to herself, (not in a very hopeful tone though), 'I won't have any pepper in my kitchen AT ALL. Soup does very well without--Maybe it's always pepper that makes people hot-tempered,'

- Sentence tokenization (aka "sentence segmentation"):
  
  But Jonathan Rees, professor of dermatology at the University of Edinburgh said it was unlikely blondes would die out completely "Genes don't die out unless there is a disadvantage of having that gene or by chance. They don't disappear," he told BBC News Online.  The only reason blondes would disappear is if having the gene was a disadvantage and I do not think that is the case.  "The frequency of blondes may drop but they won't disappear."

# Aside:
# Word segmentation and language learning

a. doyouseethekitty
b. seethedoggy
c. doyoulikethekitty
d. likethedoggy

**SEGMENTATION**

| doyou | see | thekitt | y |

| see | thedogg | y |

| doyou | like | thekitt | y |

| like | thedogg | y |

**REPRESENTATION**

LEXICON

1. doyou
2. see
3. like
4. thekitt
5. thedogg
6. y

DERIVATION

| 1 | 2 | 4 | 6 |

| 2 | 5 | 6 |

| 1 | 3 | 4 | 6 |

| 3 | 5 | 6 |

**OBJECTIVE**

**LEXICON:**
6+4+5+8+8+2 = 33

**DERIVATION:**
4+3+4+3 = 14

**TOTAL:**
33+14 = 47

# Regular Expressions

- Motivations: tokenization, morphology

- Metacharacters: .  ^  $

- Ranges: [abcdefg]  [a-g]   [^aeiou]
  \w  (word character)  [a-zA-Z0-9]
  \d   (digit)   [0-9]
  \s   (space)  [ \t\n\r\f\v]
  \W, \D, \S

- Closures: a*   a+   a?    a{3,7}

- Alternatives: (...|...|...)

- Demonstration (nltk.app.nemo)