

[illegible]

Introducing the tasks:  
Getting simple structured  
information out of text



# Information Extraction

- Information extraction (IE) systems
  - Find and understand limited relevant parts of texts
  - Gather information from many pieces of text
  - Produce a structured representation of relevant information:
    - *relations* (in the database sense), a.k.a.,
    - *a knowledge base*
  - Goals:
    1. Organize information so that it is useful to people
    2. Put information in a semantically precise form that allows further inferences to be made by computer algorithms



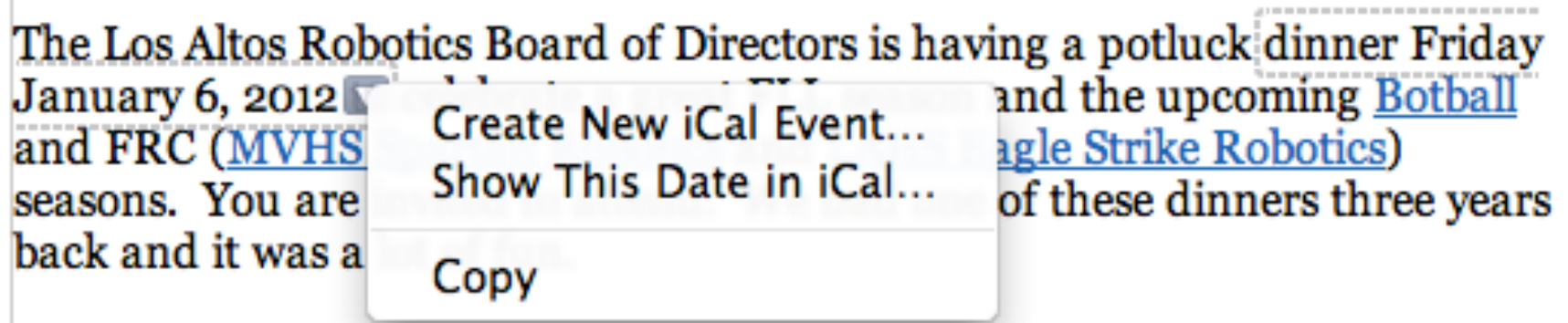
# Information Extraction (IE)

- IE systems extract clear, factual information
  - Roughly: *Who did what to whom when?*
- E.g.,
  - Gathering earnings, profits, board members, headquarters, etc. from company reports
    - The headquarters of BHP Billiton Limited, and the global headquarters of the combined BHP Billiton Group, are located in Melbourne, Australia.
    - *headquarters("BHP Biliton Limited", "Melbourne, Australia")*
  - Learn drug-gene product interactions from medical research literature



## Low-level information extraction

- Is now available – and I think popular – in applications like Apple or Google mail, and web indexing



- Often seems to be based on regular expressions and name lists



# Low-level information extraction



bhp billiton headquarters

Search

About 123,000 results (0.23 seconds)

Everything

Best guess for BHP Billiton Ltd. Headquarters is **Melbourne, London**

Images

Mentioned on at least 9 websites including [wikipedia.org](http://wikipedia.org), [bhpbilliton.com](http://bhpbilliton.com) and [bhpbilliton.com](http://bhpbilliton.com) - [Feedback](#)

Maps

[BHP Billiton - Wikipedia, the free encyclopedia](http://en.wikipedia.org/wiki/BHP_Billiton)

Videos

[en.wikipedia.org/wiki/BHP\\_Billiton](http://en.wikipedia.org/wiki/BHP_Billiton)

News

Merger of BHP & Billiton 2001 (creation of a DLC). **Headquarters, Melbourne, Australia (BHP Billiton Limited and BHP Billiton Group) London, United Kingdom ...**

Shopping

[History](#) - [Corporate affairs](#) - [Operations](#) - [Accidents](#)





# Named Entity Recognition (NER)

- A very important sub-task: **find** and **classify** names in text, for example:
  - The decision by the independent MP Andrew Wilkie to withdraw his support for the minority Labor government sounded dramatic but it should not further threaten its stability. When, after the 2010 election, Wilkie, Rob Oakeshott, Tony Windsor and the Greens agreed to support Labor, they gave just two guarantees: confidence and supply.



# Named Entity Recognition (NER)

- A very important sub-task: **find** and **classify** names in text, for example:
  - The decision by the independent MP **Andrew Wilkie** to withdraw his support for the minority **Labor** government sounded dramatic but it should not further threaten its stability. When, after the **2010** election, **Wilkie**, **Rob Oakeshott**, **Tony Windsor** and the **Greens** agreed to support **Labor**, they gave just two guarantees: confidence and supply.



# Named Entity Recognition (NER)

- A very important sub-task: **find** and **classify** names in text, for example:
  - The decision by the independent MP **Andrew Wilkie** to withdraw his support for the minority **Labor** government sounded dramatic but it should not further threaten its stability. When, after the **2010** election, **Wilkie**, **Rob Oakeshott**, **Tony Windsor** and the **Greens** agreed to support **Labor**, they gave just two guarantees: confidence and supply.

<b>Person</b>
<b>Date</b>
<b>Location</b>
<b>Organi- zation</b>

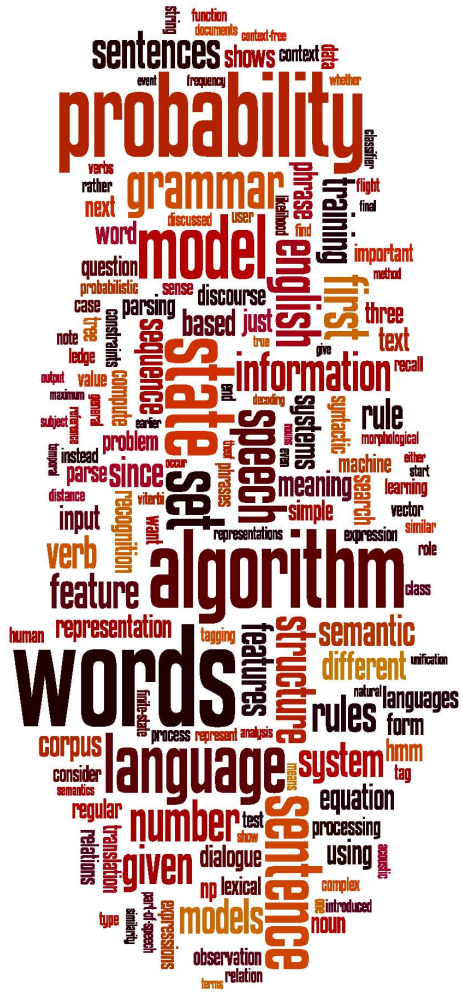




# Named Entity Recognition (NER)

- The uses:
  - Named entities can be indexed, linked off, etc.
  - Sentiment can be attributed to companies or products
  - A lot of IE relations are associations between named entities
  - For question answering, answers are often named entities.
- Concretely:
  - Many web pages tag various entities, with links to bio or topic pages, etc.
    - Reuters' OpenCalais, Evri, AlchemyAPI, Yahoo's Term Extraction, ...
  - Apple/Google/Microsoft/... smart recognizers for document content

Introducing the tasks:  
Getting simple structured  
information out of text



# Evaluation of Named Entity Recognition

## The extension of Precision, Recall, and the F measure to sequences



# The Named Entity Recognition Task

Task: Predict entities in a text

Foreign      **ORG**

Ministry      **ORG**

spokesman      **O**

Shen      **PER**

Guofang      **PER**

told      **O**

Reuters      **ORG**

:

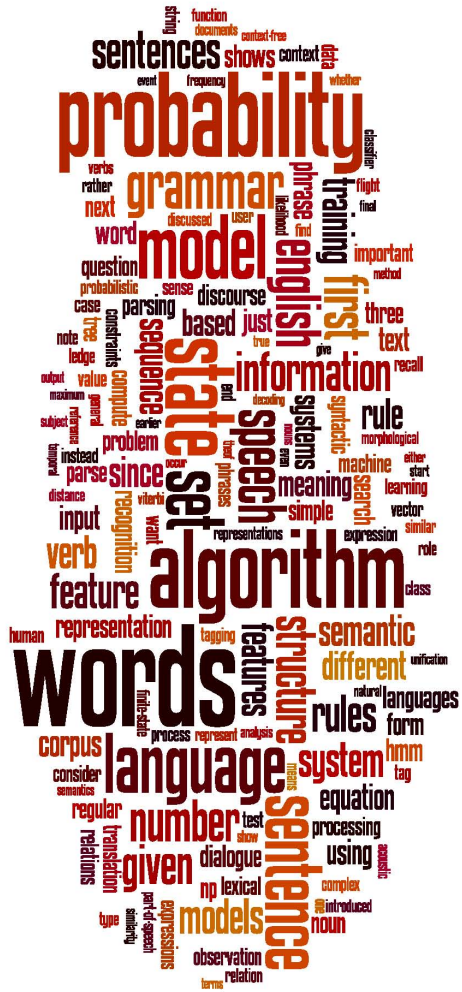


Standard  
evaluation  
is per entity,  
*not* per token



## Precision/Recall/F1 for IE/NER

- Recall and precision are straightforward for tasks like IR and text categorization, where there is only one grain size (documents)
- The measure behaves a bit funnily for IE/NER when there are *boundary errors* (which are *common*):
  - First Bank of Chicago announced earnings ...
- This counts as both a fp and a fn
- Selecting *nothing* would have been better
- Some other metrics (e.g., MUC scorer) give partial credit (according to complex rules)



# Evaluation of Named Entity Recognition

## The extension of Precision, Recall, and the F measure to sequences



# Sequence Models for Named Entity Recognition



# The ML sequence model approach to NER

## Training

1. Collect a set of representative training documents
2. Label each token for its entity class or other (O)
3. Design feature extractors appropriate to the text and classes
4. Train a sequence classifier to predict the labels from the data

## Testing

1. Receive a set of testing documents
2. Run sequence model inference to label each token
3. Appropriately output the recognized entities





## Encoding classes for sequence labeling

	IO encoding	IOB encoding
Fred	PER	B-PER
showed	O	O
Sue	PER	B-PER
Mengqiu	PER	B-PER
Huang	PER	I-PER
's	O	O
new	O	O
painting	O	O

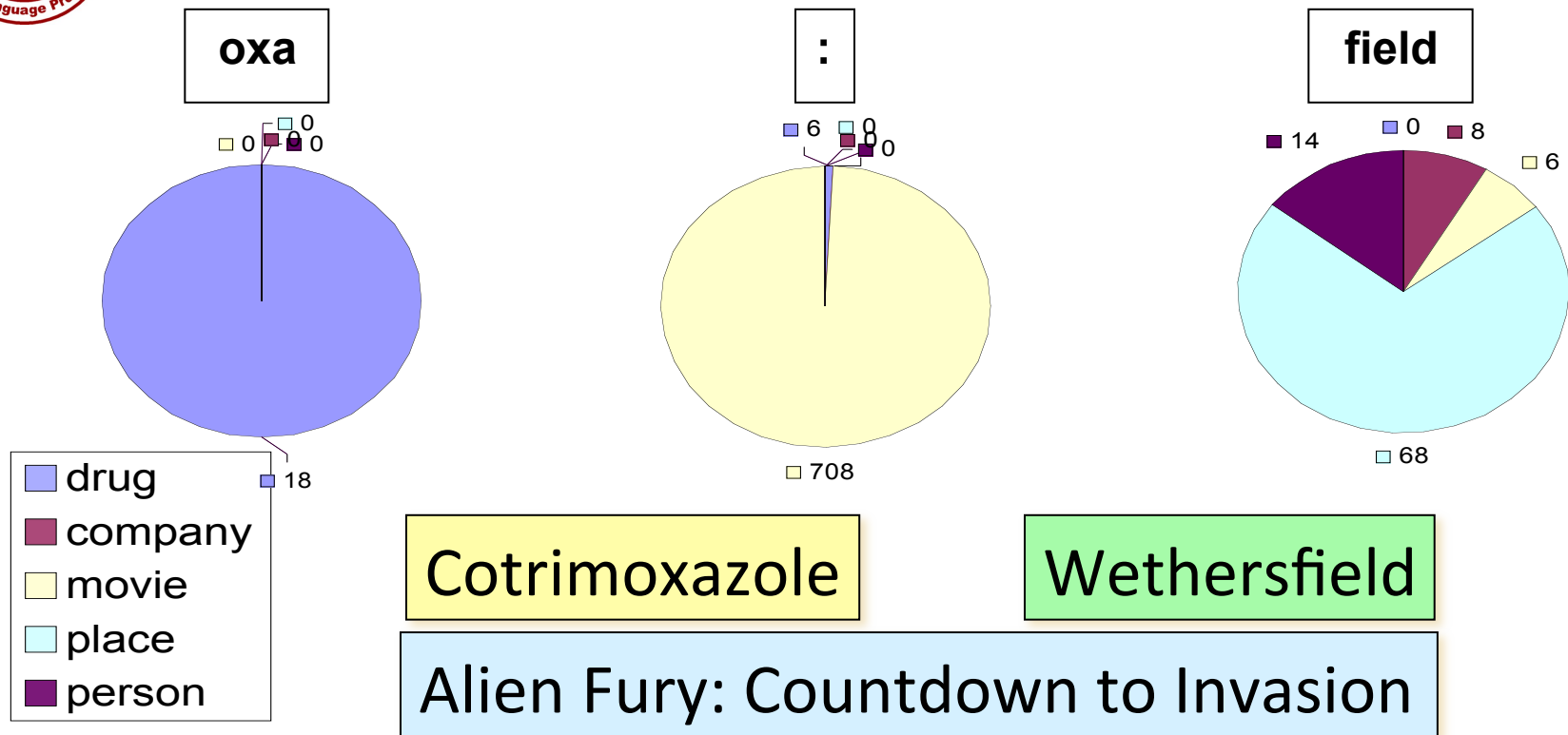


# Features for sequence labeling

- Words
  - Current word (essentially like a learned dictionary)
  - Previous/next word (context)
- Other kinds of inferred linguistic classification
  - Part-of-speech tags
- Label context
  - Previous (and perhaps next) label



## Features: Word substrings

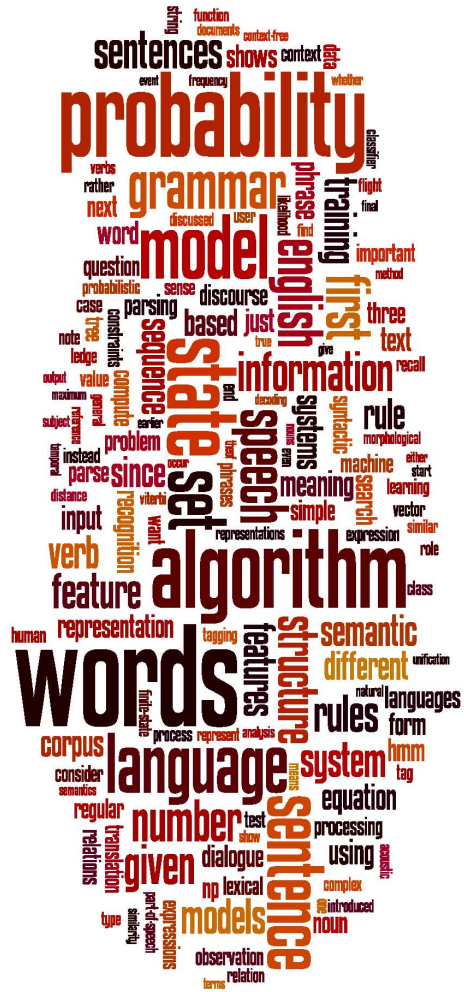




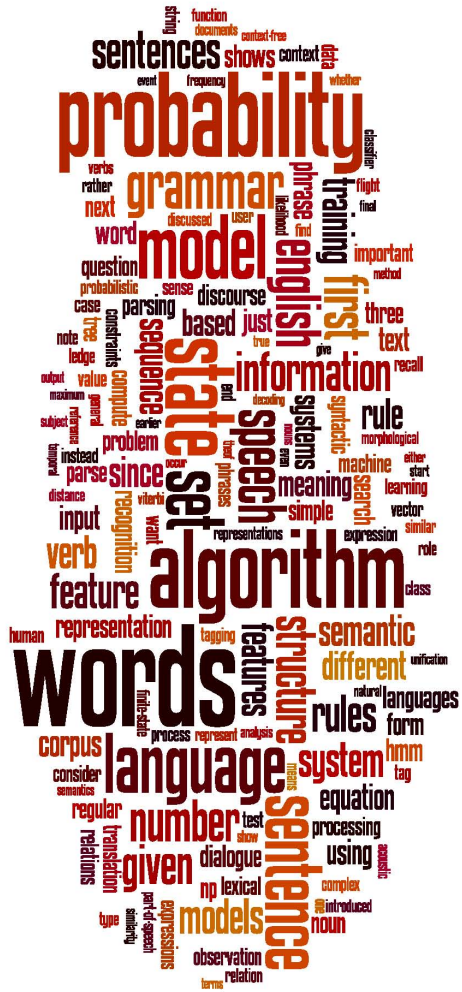
# Features: Word shapes

- Word Shapes
  - Map words to simplified representation that encodes attributes such as length, capitalization, numerals, Greek letters, internal punctuation, etc.

Varicella-zoster	Xx-xxx
mRNA	xXXX
CPA1	XXXd



# Sequence Models for Named Entity Recognition



# Maximum entropy sequence models

Maximum entropy Markov models (MEMMs) or Conditional Markov models



# Sequence problems

- Many problems in NLP have data which is a sequence of characters, words, phrases, lines, or sentences ...
- We can think of our task as one of labeling each item

VBG	NN	IN	DT	NN	IN	NN
Chasing	opportunity	in	an	age	of	upheaval

## POS tagging

PERS	O	O	O	ORG	ORG
Murdoch	discusses	future	of	News	Corp.

## Named entity recognition

B	B	I	I	B	I	B	I	B	B
而	相	对	于	这	些	品	牌	的	价

## Word segmentation

Q
A
Q
A
A
A
Q
A

## Text segmentation



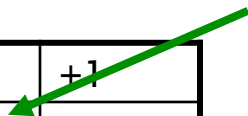
# MEMM inference in systems

- For a Conditional Markov Model (CMM) a.k.a. a Maximum Entropy Markov Model (MEMM), the classifier makes a single decision at a time, conditioned on evidence from observations **and previous decisions**
- A larger space of sequences is usually explored via search

Local Context

-3	-2	-1	0	+1
DT	NNP	VBD	???	???
The	Dow	fell	22.6	%

Decision Point



Features

$W_0$	22.6
$W_{+1}$	%
$W_{-1}$	fell
$T_{-1}$	VBD
$T_{-1}-T_{-2}$	NNP-VBD
hasDigit?	true
...	...

(Ratnaparkhi 1996; Toutanova et al. 2003, etc.)





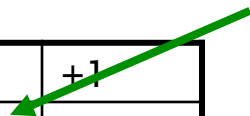
## Example: POS Tagging

- Scoring individual labeling decisions is no more complex than standard classification decisions
  - We have some assumed labels to use for prior positions
  - We use features of those and the observed data (which can include current, previous, and next words) to predict the current label

### Local Context

-3	-2	-1	0	+1
DT	NNP	VBD	???	???
The	Dow	fell	22.6	%

### Decision Point



### Features

$W_0$	22.6
$W_{+1}$	%
$W_{-1}$	fell
$T_{-1}$	VBD
$T_{-1}-T_{-2}$	NNP-VBD
hasDigit?	true
...	...

(Ratnaparkhi 1996; Toutanova et al. 2003, etc.)



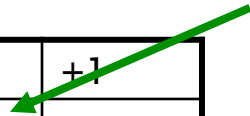
## Example: POS Tagging

- POS tagging Features can include:
  - Current, previous, next words in isolation or together.
  - Previous one, two, three tags.
  - Word-internal features: word types, suffixes, dashes, etc.

### Local Context

-3	-2	-1	0	+1
DT	NNP	VBD	???	???
The	Dow	fell	22.6	%

### Decision Point



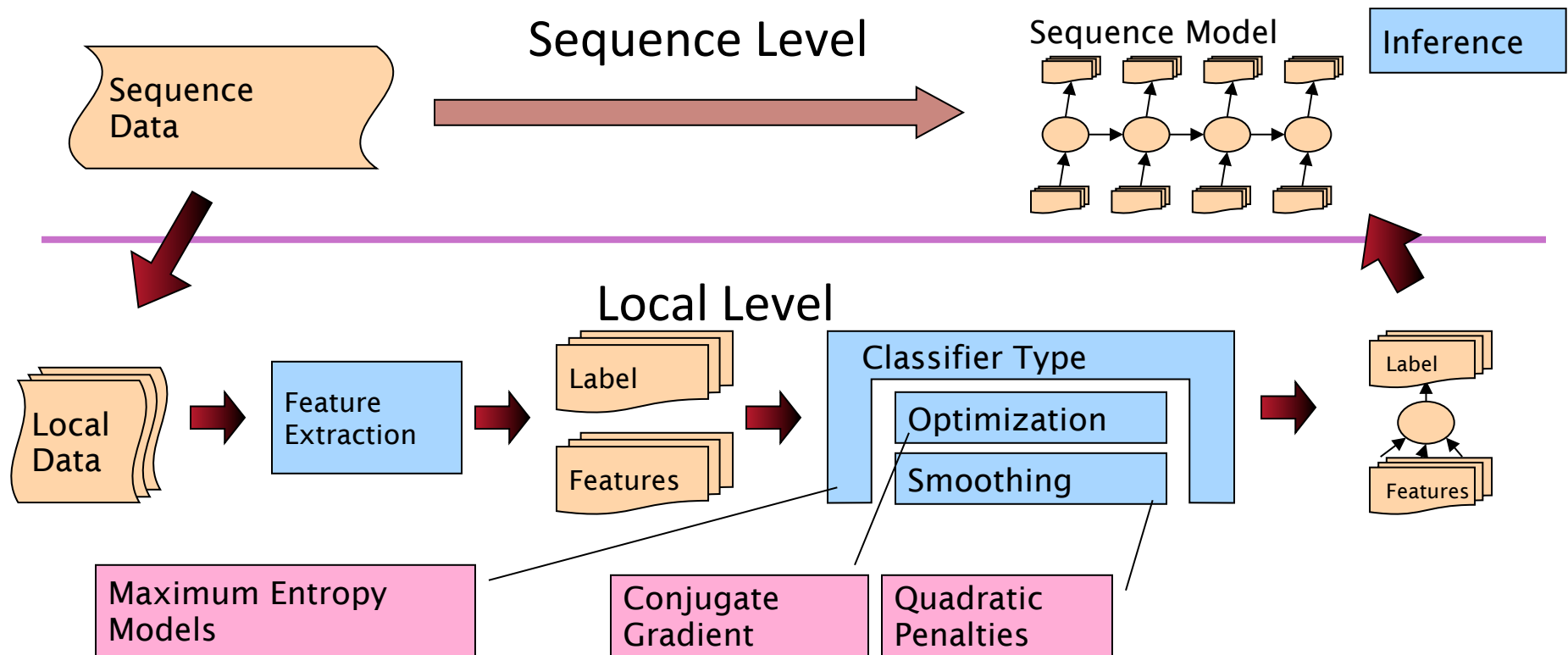
### Features

$W_0$	22.6
$W_{+1}$	%
$W_{-1}$	fell
$T_{-1}$	VBD
$T_{-1}-T_{-2}$	NNP-VBD
hasDigit?	true
...	...

(Ratnaparkhi 1996; Toutanova et al. 2003, etc.)

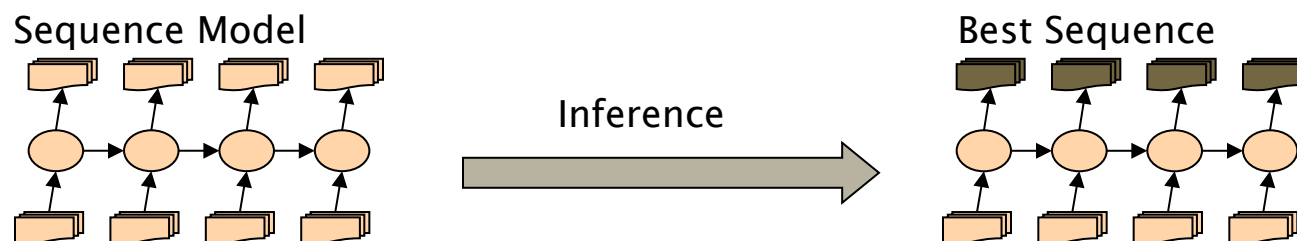


# Inference in Systems





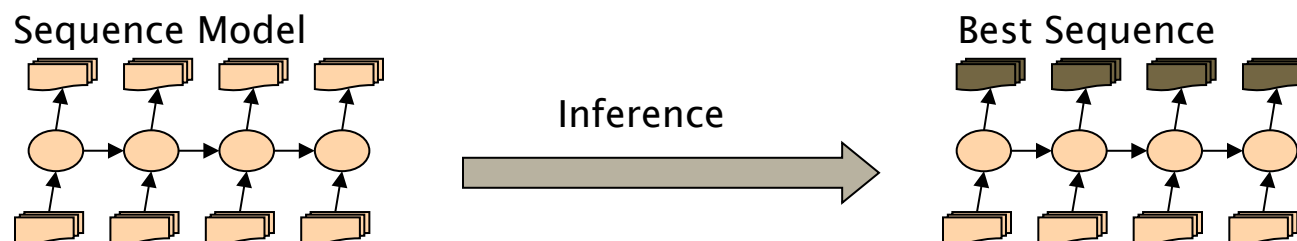
# Greedy Inference



- Greedy inference:
  - We just start at the left, and use our classifier at each position to assign a label
  - The classifier can depend on previous labeling decisions as well as observed data
- Advantages:
  - Fast, no extra memory requirements
  - Very easy to implement
  - With rich features including observations to the right, it may perform quite well
- Disadvantage:
  - Greedy. We make commit errors we cannot recover from



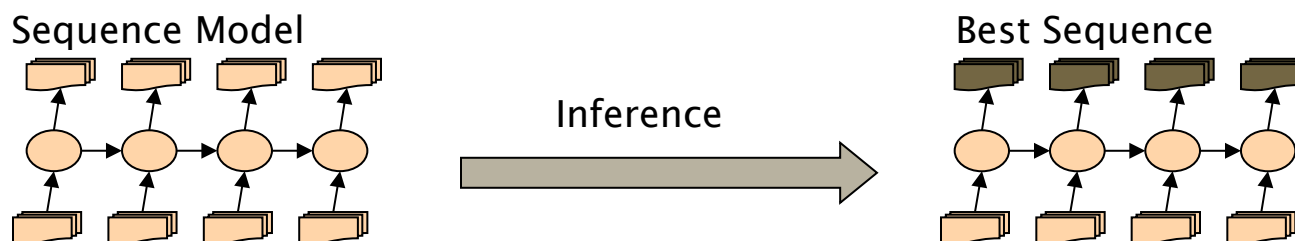
# Beam Inference



- Beam inference:
  - At each position keep the top  $k$  complete sequences.
  - Extend each sequence in each local way.
  - The extensions compete for the  $k$  slots at the next position.
- Advantages:
  - Fast; beam sizes of 3–5 are almost as good as exact inference in many cases.
  - Easy to implement (no dynamic programming required).
- Disadvantage:
  - Inexact: the globally best sequence can fall off the beam.



# Viterbi Inference



- Viterbi inference:
  - Dynamic programming or memoization.
  - Requires small window of state influence (e.g., past two states are relevant).
- Advantage:
  - Exact: the global best sequence is returned.
- Disadvantage:
  - Harder to implement long-distance state-state interactions (but beam inference tends not to allow long-distance resurrection of sequences anyway).



## CRFs [Lafferty, Pereira, and McCallum 2001]

- Another sequence model: Conditional Random Fields (CRFs)
- A whole-sequence conditional model rather than a chaining of local models.

$$P(c \mid d, \lambda) = \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}$$

- The space of  $c$ 's is now the space of sequences
  - But if the features  $f_i$  remain local, the conditional sequence likelihood can be calculated exactly using dynamic programming
- Training is slower, but CRFs avoid causal-competition biases
- These (or a variant using a max margin criterion) are seen as the state-of-the-art these days ... but in practice usually work much the same as MEMMs.



# Maximum entropy sequence models

# Maximum entropy Markov models (MEMMs) or Conditional Markov models