

How to Read *War and Peace* in 30 Seconds

Denis Griffis

PhD Student, Speech and Language Technologies lab
The Ohio State University



THE OHIO STATE
UNIVERSITY

COLLEGE OF ENGINEERING

language*speech*und... 0101010010010101
speech & 0100101001010101
language 0101010101010101
technologies 00111110000101
@osu 001010010101
recognition*dialogue*gen... 010101010101

How to Read *War and Peace* in 30 Seconds

Or: An introduction to Natural Language Processing

Denis Griffis

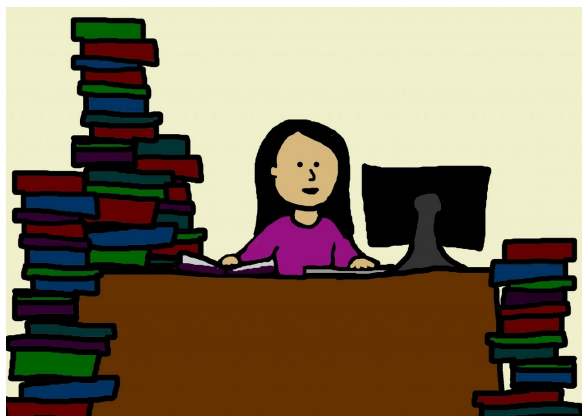
PhD Student, Speech and Language Technologies lab
The Ohio State University

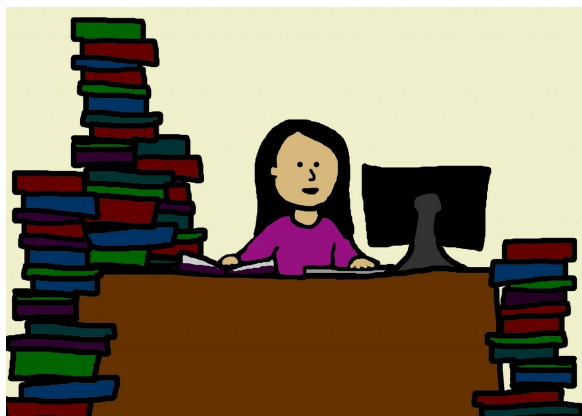


THE OHIO STATE
UNIVERSITY

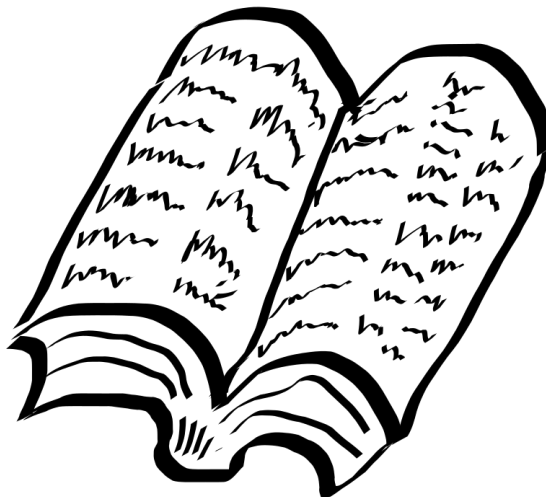
COLLEGE OF ENGINEERING



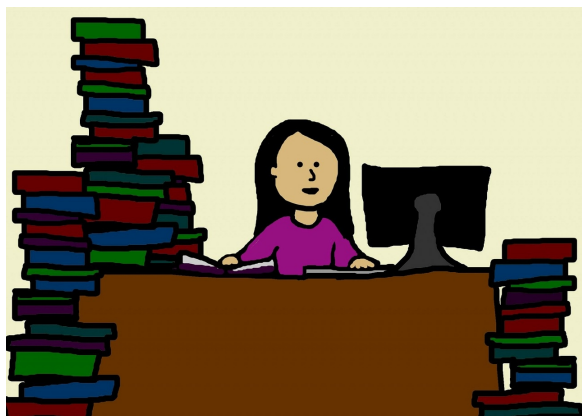




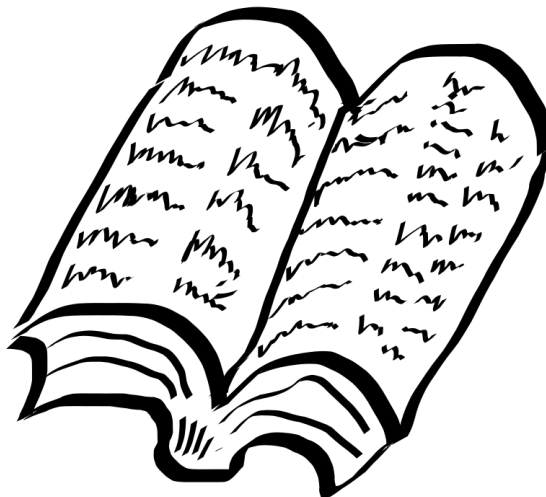
+



=

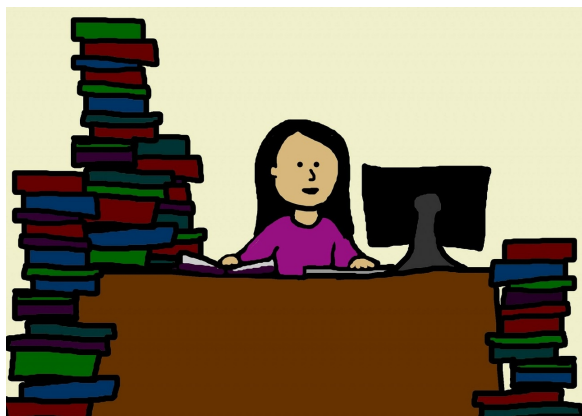


+



=



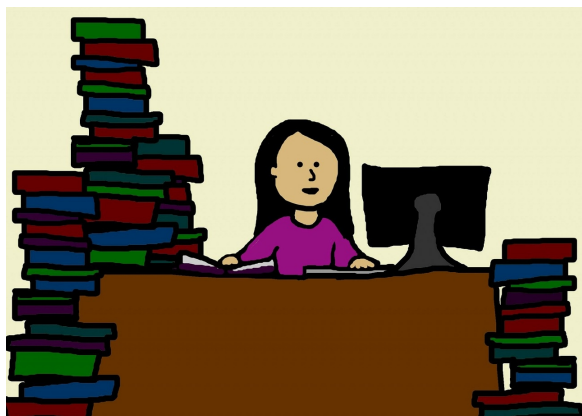


+

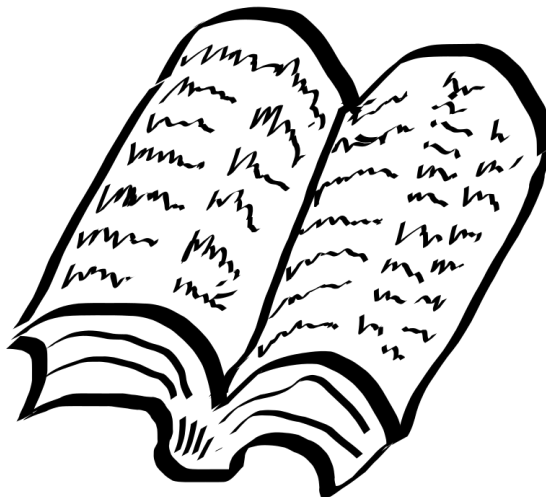


=





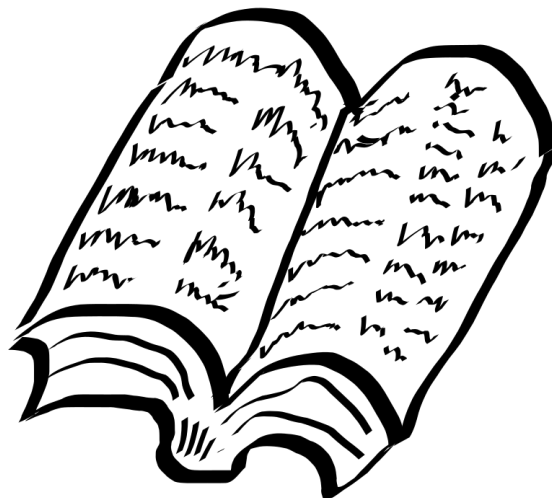
+



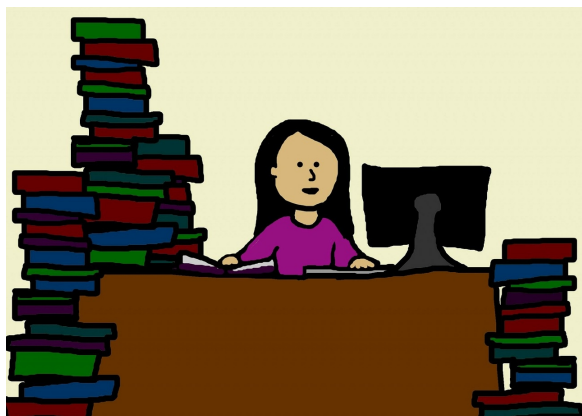
=



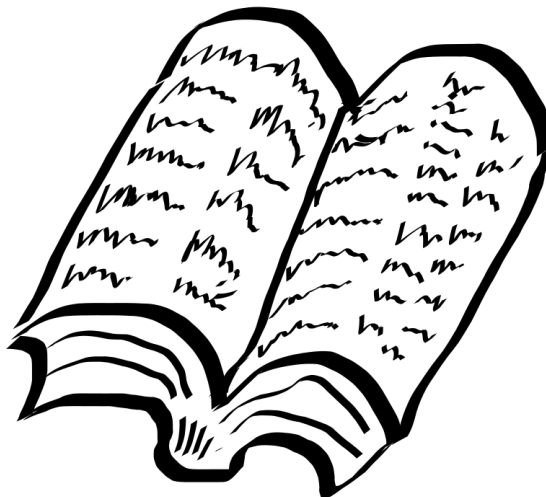
+



=



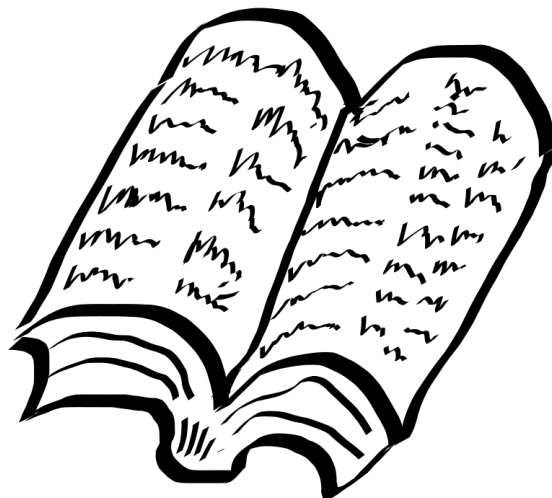
+



=



+



=





How the heck do we
get a computer to
understand text?

Natural language processing

From Wikipedia, the free encyclopedia

This article is about language processing by computers. For the processing of language by the human brain, see [Language processing](#).

Natural language processing (NLP) is a field of [computer science](#), [artificial intelligence](#), and [computational linguistics](#) concerned with the interactions between computers and human (natural) languages. As such, NLP is related to the area of [human–computer interaction](#). Many





OK Google,
where are all
the cats?



OK Google,
where are all
the cats?



```
SELECT CurrentLocation  
FROM AllAnimals  
WHERE AnimalType =  
'Cat'
```

Speech Recognition



OK Google,
where are all
the cats?



```
SELECT CurrentLocation  
FROM AllAnimals  
WHERE AnimalType =  
'Cat'
```

Speech Recognition



OK Google,
where are all
the cats?



```
SELECT CurrentLocation  
FROM AllAnimals  
WHERE AnimalType =  
'Cat'
```

Natural Language Processing

“Yess! Yess! Its official Nintendo announced today that they Will release the Nintendo 3DS in north America march 27 for \$250”

*“Yess! Yess! Its official [**Nintendo**] announced today that they Will release the [**Nintendo 3DS**] in [**north America**] [**march 27**] for [**\$250**]”*

*“Yess! Yess! Its official [**Nintendo**] announced today that they Will release the [**Nintendo 3DS**] in [**north America**] [**march 27**] for [**\$250**]”*



| <i>Company</i> | <i>Product</i> | <i>Date</i> | <i>Price</i> | <i>Region</i> |
|----------------|----------------|-------------|--------------|---------------|
| Nintendo | Nintendo 3DS | March 27 | \$250 | North America |

*“Yess! Yess! Its official [**Nintendo**] announced today that they Will release the [**Nintendo 3DS**] in [**north America**] [**march 27**] for [**\$250**]”*



| <i>Company</i> | <i>Product</i> | <i>Date</i> | <i>Price</i> | <i>Region</i> |
|----------------|----------------|-------------|--------------|---------------|
| Nintendo | Nintendo 3DS | March 27 | \$250 | North America |

Natural language *understanding*

| <i>Company</i> | <i>Product</i> | <i>Date</i> | <i>Price</i> | <i>Region</i> |
|----------------|----------------|-------------|--------------|---------------|
| Nintendo | Nintendo 3DS | March 27 | \$250 | North America |

*Nintendo will release the Nintendo 3DS
in North America on March 27 for \$250.*



| <i>Company</i> | <i>Product</i> | <i>Date</i> | <i>Price</i> | <i>Region</i> |
|----------------|----------------|-------------|--------------|---------------|
| Nintendo | Nintendo 3DS | March 27 | \$250 | North America |

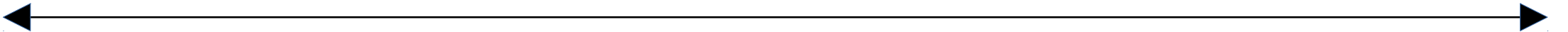
*Nintendo will release the Nintendo 3DS
in North America on March 27 for \$250.*



| <i>Company</i> | <i>Product</i> | <i>Date</i> | <i>Price</i> | <i>Region</i> |
|----------------|----------------|-------------|--------------|---------------|
| Nintendo | Nintendo 3DS | March 27 | \$250 | North America |

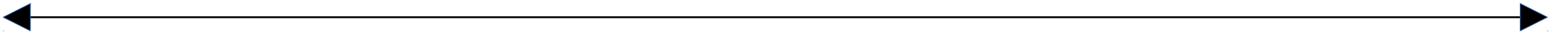
Natural language *generation*

A Brief History of NLP



A Brief History of NLP

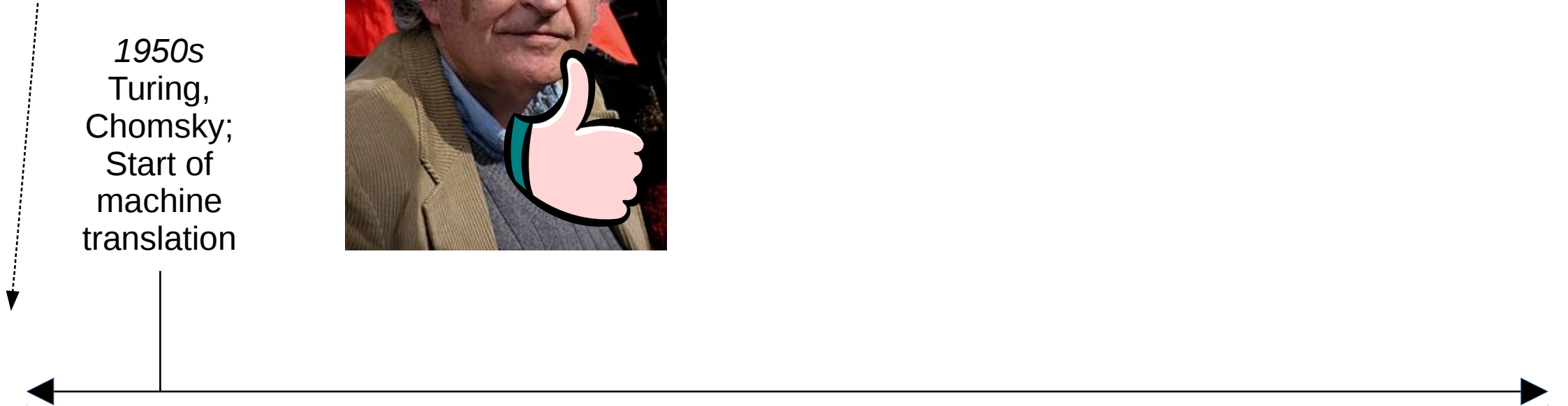
< 1950
Codes for
translation;
message
encodings



A Brief History of NLP

< 1950
Codes for
translation;
message
encodings

1950s
Turing,
Chomsky;
Start of
machine
translation

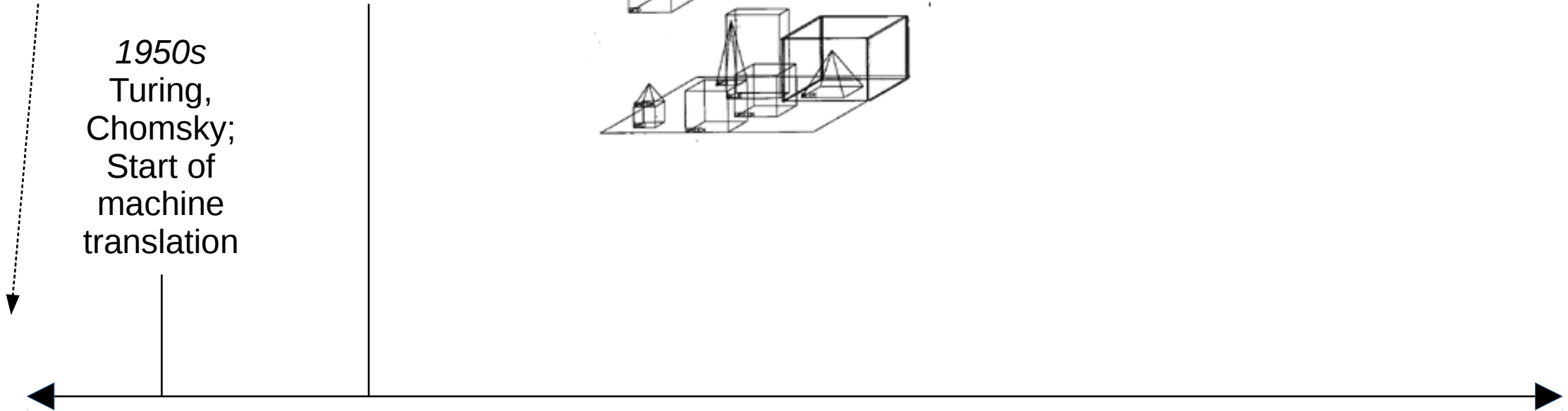
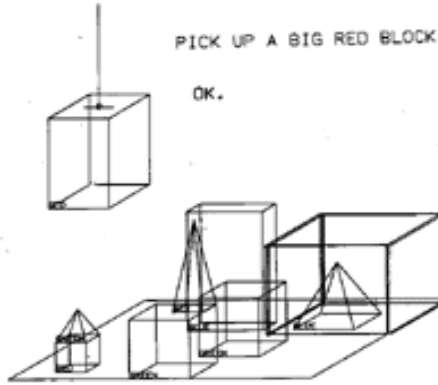


A Brief History of NLP

< 1950
Codes for
translation;
message
encodings

Late 1960s
SHRDLU,
ELIZA

1950s
Turing,
Chomsky;
Start of
machine
translation



A Brief History of NLP

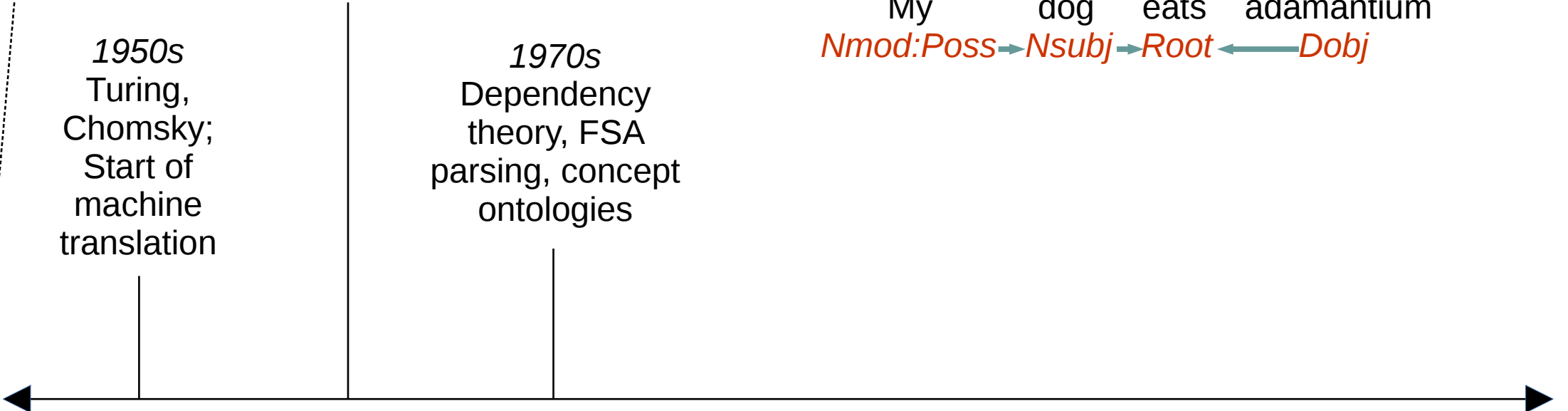
< 1950
Codes for
translation;
message
encodings

Late 1960s
SHRDLU,
ELIZA

1950s
Turing,
Chomsky;
Start of
machine
translation

1970s
Dependency
theory, FSA
parsing, concept
ontologies

My dog eats adamantium
Nmod:Poss → *Nsubj* → *Root* ← *Dobj*



A Brief History of NLP

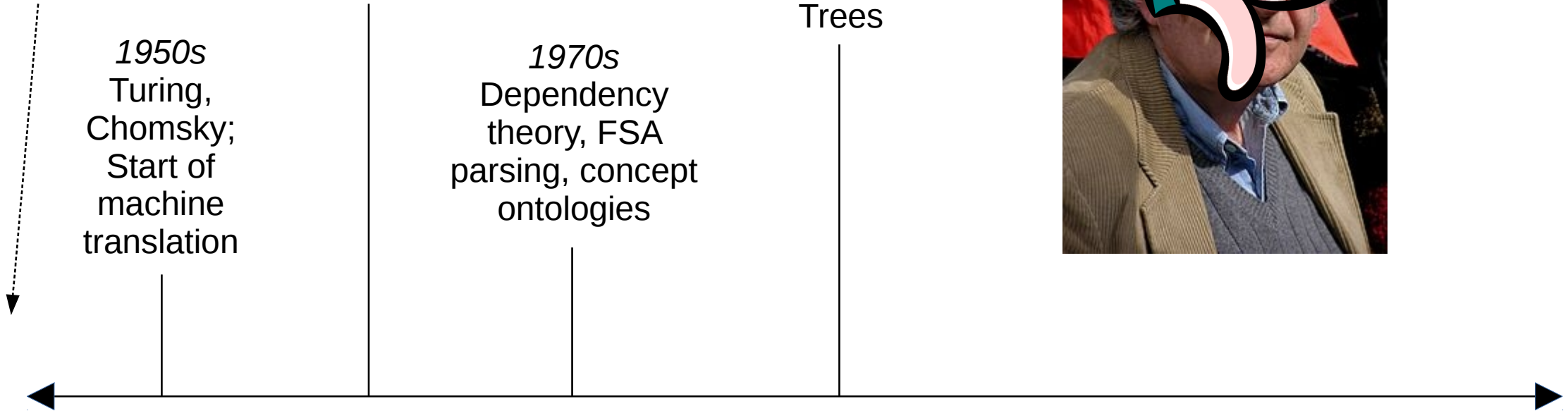
< 1950
Codes for
translation;
message
encodings

Late 1960s
SHRDLU,
ELIZA

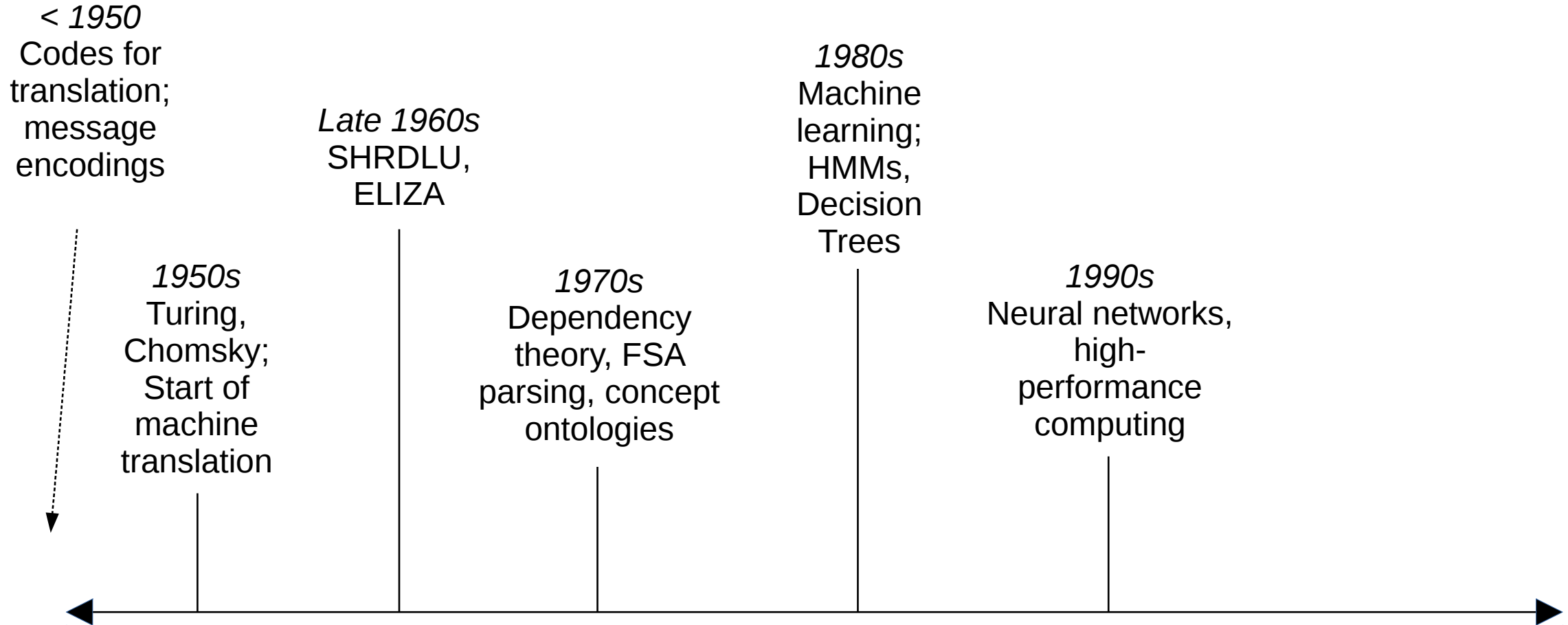
1980s
Machine
learning;
HMMs,
Decision
Trees

1950s
Turing,
Chomsky;
Start of
machine
translation

1970s
Dependency
theory, FSA
parsing, concept
ontologies



A Brief History of NLP



A Brief History of NLP



< 1950
Codes for
translation;
message
encodings

Late 1960s
SHRDLU,
ELIZA

1980s
Machine
learning;
HMMs,
Decision
Trees

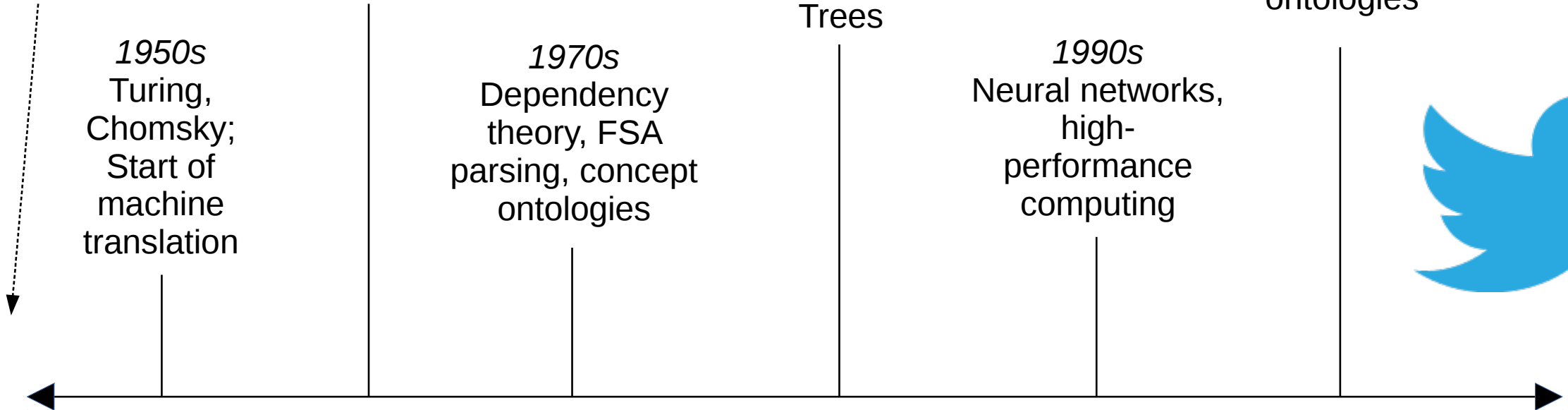


2000s on
Explosion in Web
text data,
ontologies

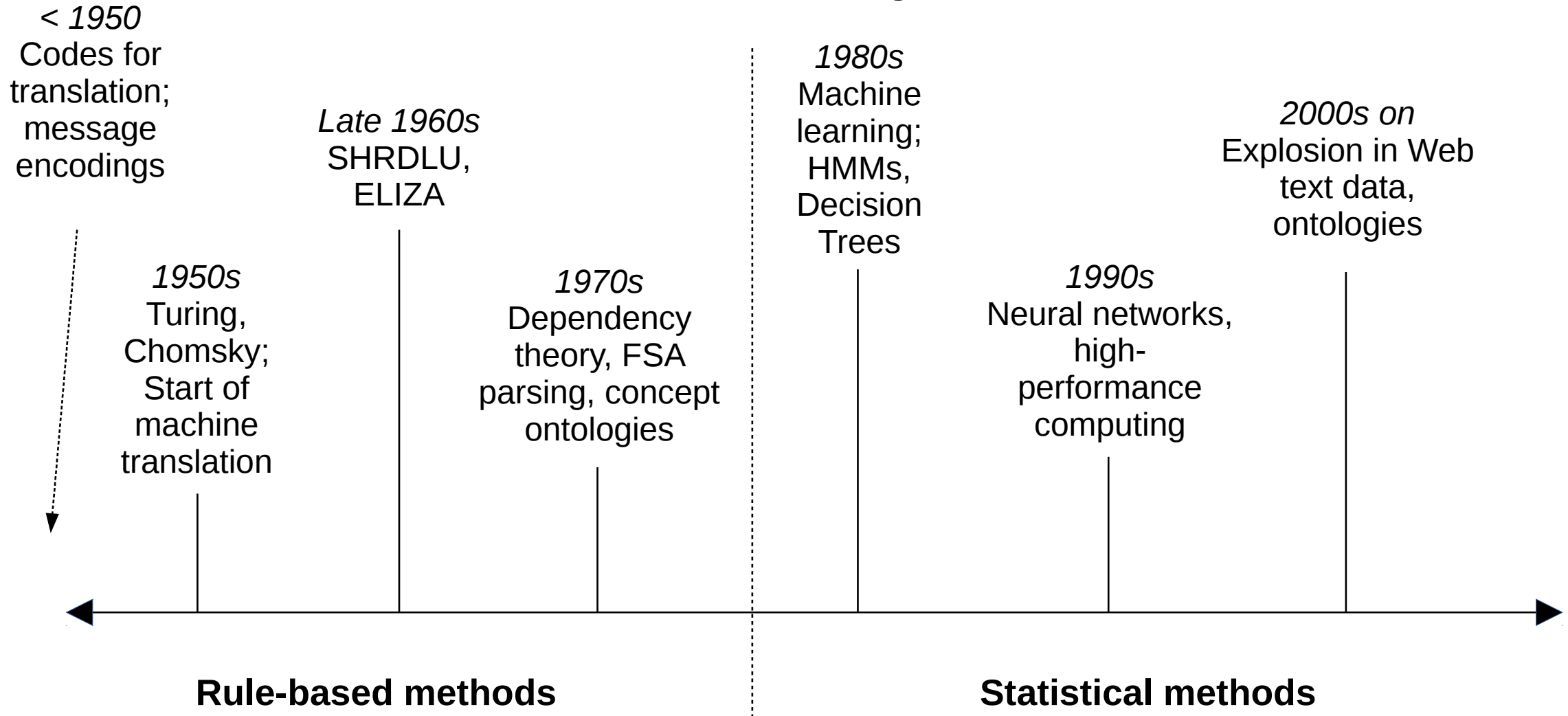
1950s
Turing,
Chomsky;
Start of
machine
translation

1970s
Dependency
theory, FSA
parsing, concept
ontologies

1990s
Neural networks,
high-
performance
computing



A Brief History of NLP



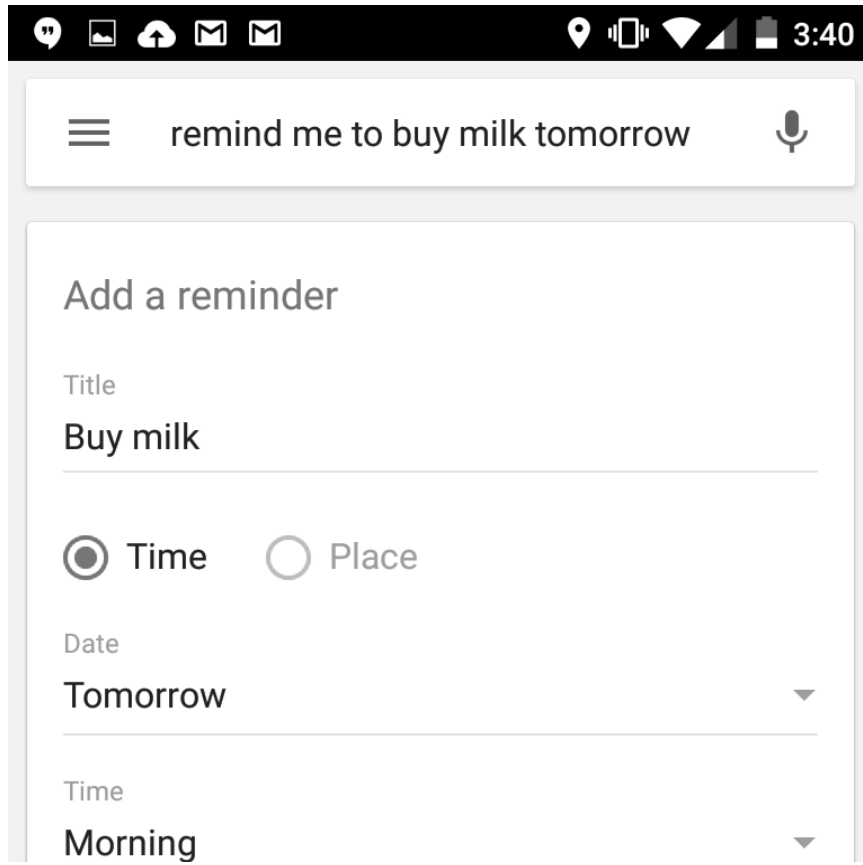
Rule-Based NLP

- Regular Expressions

```
#####  
#      Per cause_of_death  
#####  
  
{  
    ruletype: "composite",  
    pattern: (([ner:PERSON]+) /died/ /of|from/ /a/? ([tag:NN]+)),  
    result: Format("per:cause_of_death(%s,%s)", $1.word, $2.word),  
    action: (Annotate($1, kbp, "per"), Annotate($2, kbp, "per_cause_of_death"))  
}
```

Rule-Based NLP

- Keywords and arguments



remind me to buy milk tomorrow

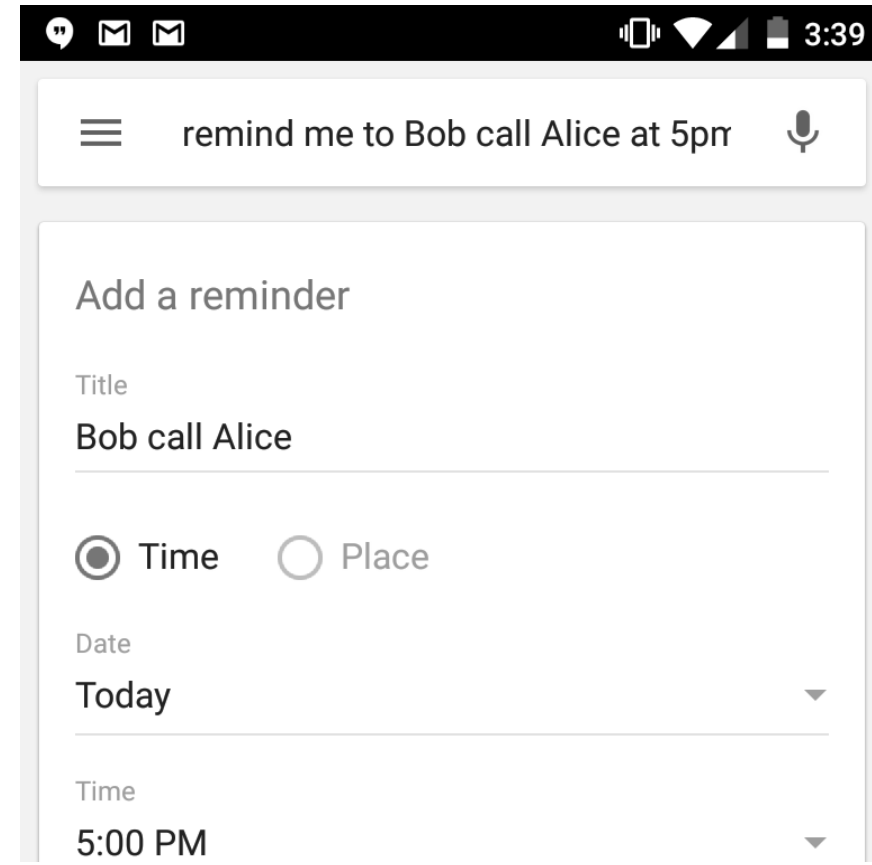
Add a reminder

Title
Buy milk

Time Place

Date
Tomorrow

Time
Morning



remind me to Bob call Alice at 5pr

Add a reminder

Title
Bob call Alice

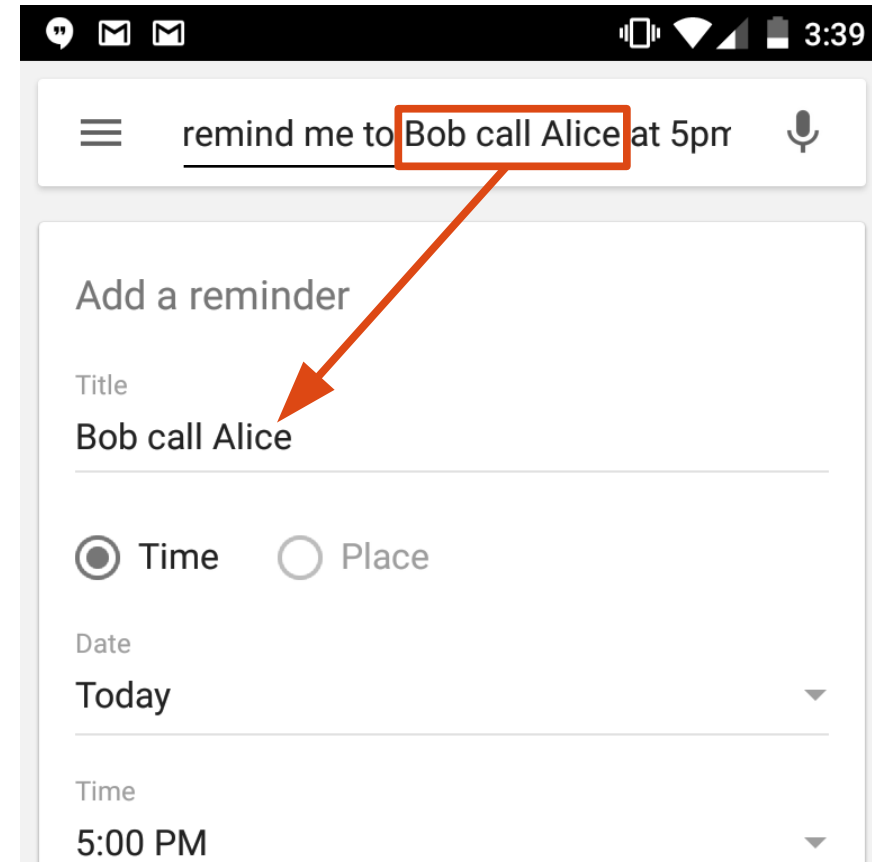
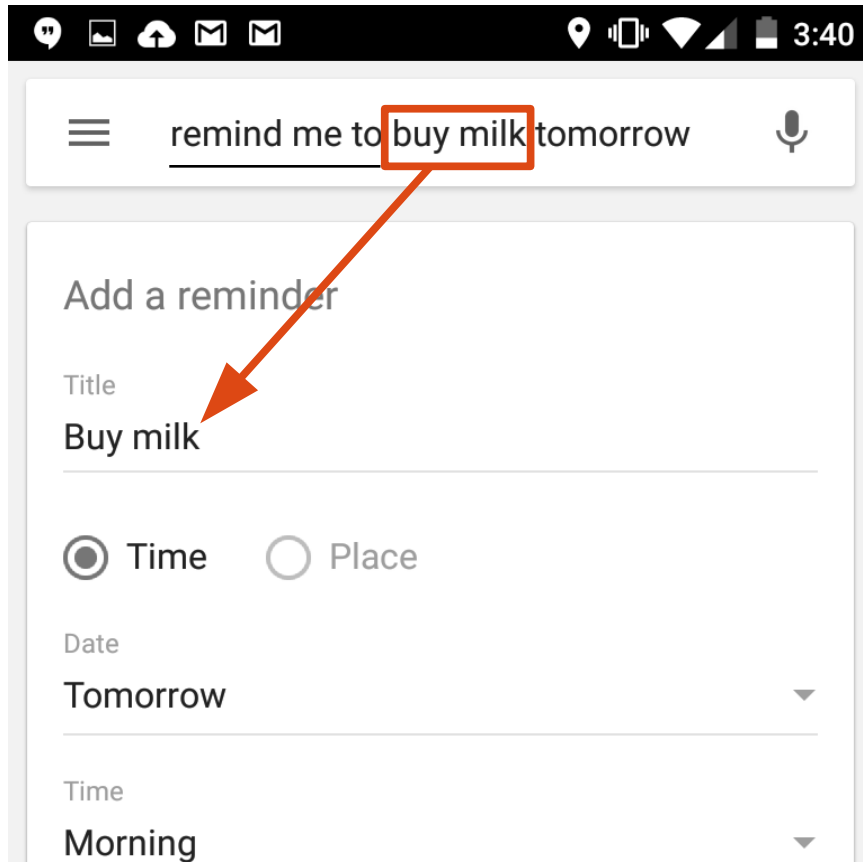
Time Place

Date
Today

Time
5:00 PM

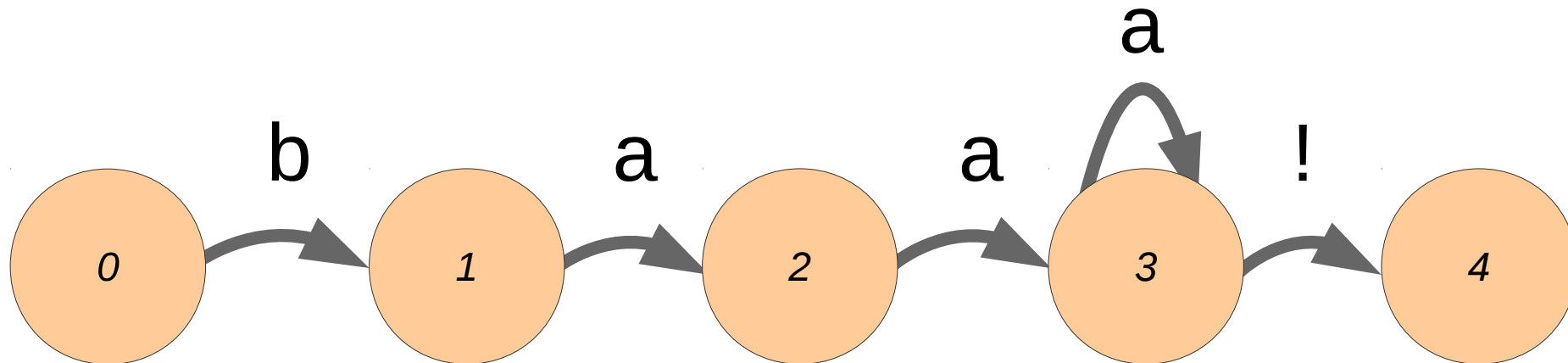
Rule-Based NLP

- Keywords and arguments



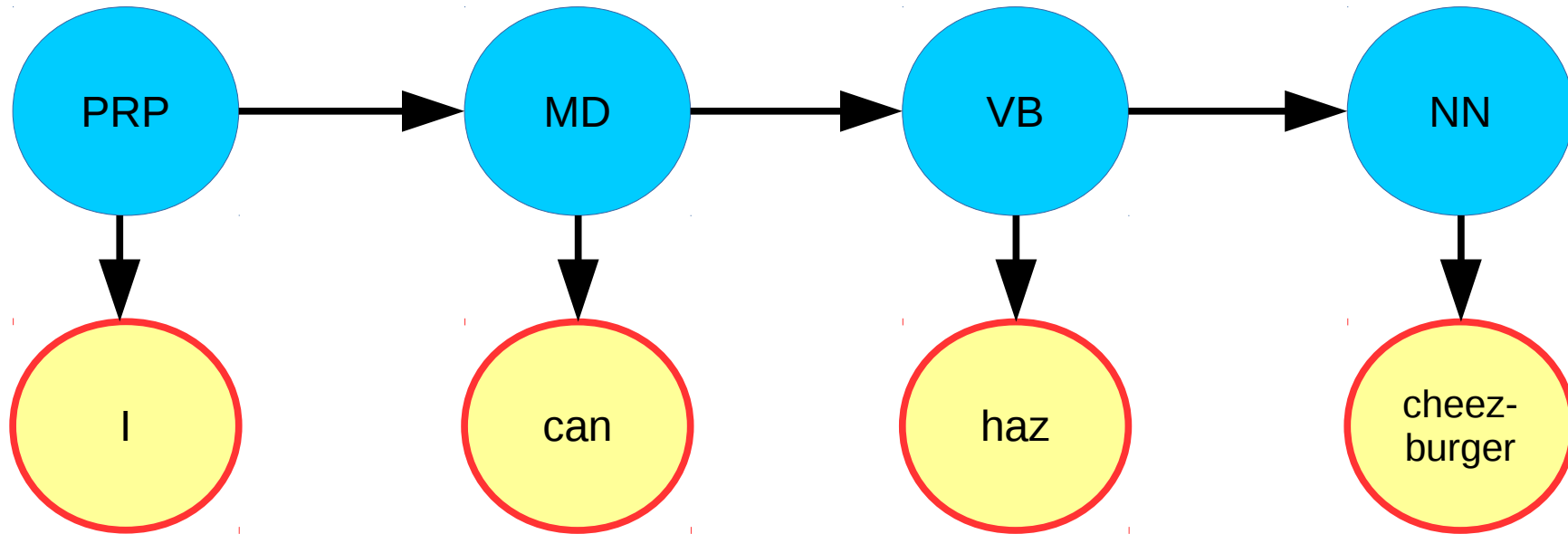
Rule-Based NLP

- Finite State Automata



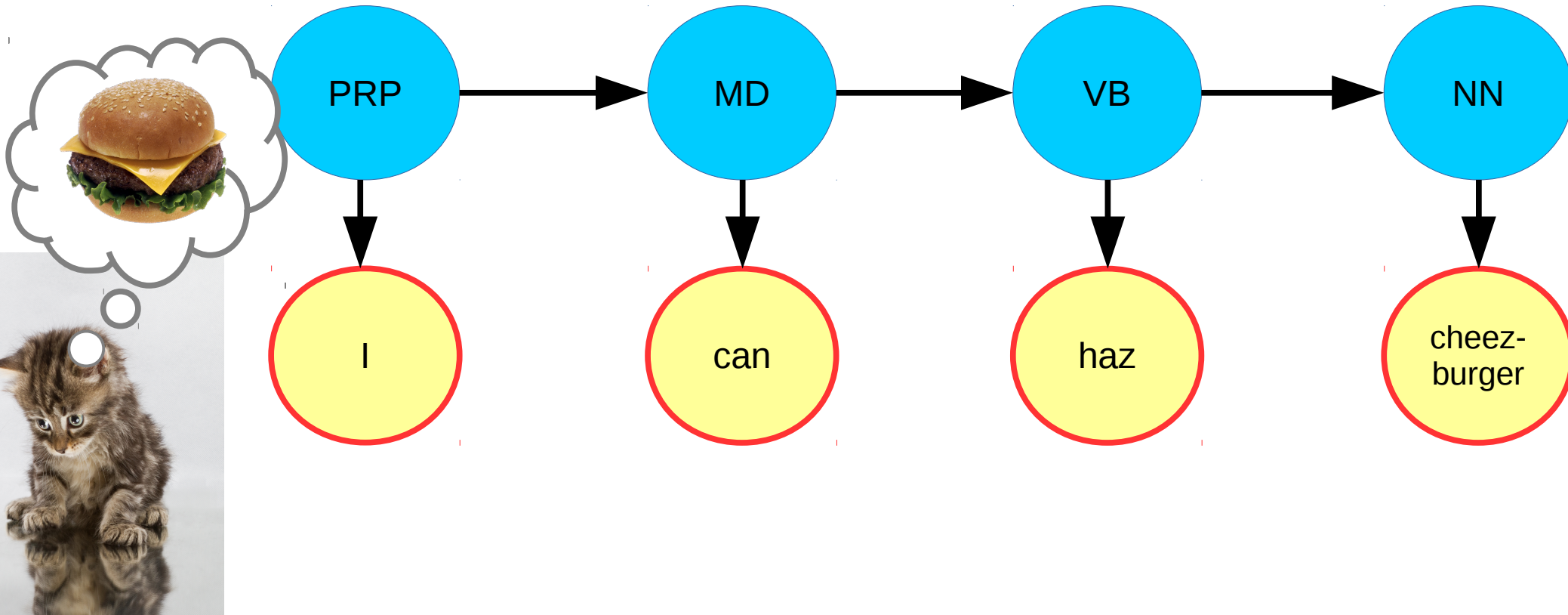
Statistical NLP

- Hidden Markov Models (HMMs)



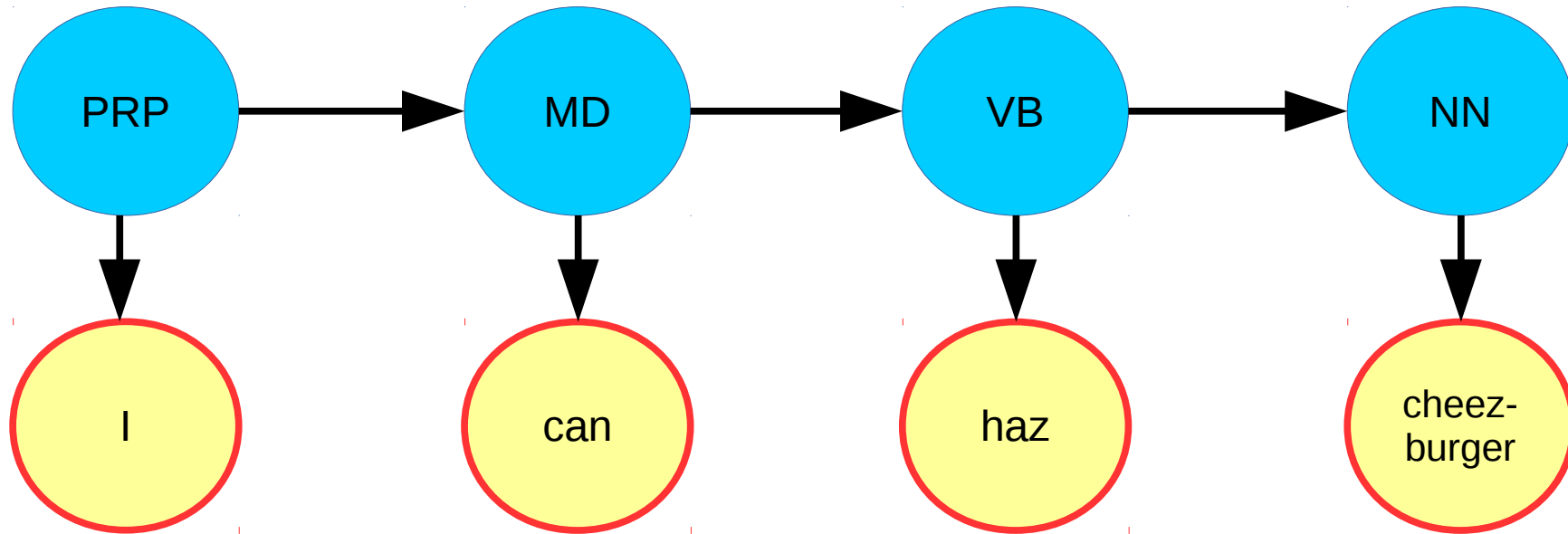
Statistical NLP

- Hidden Markov Models (HMMs)



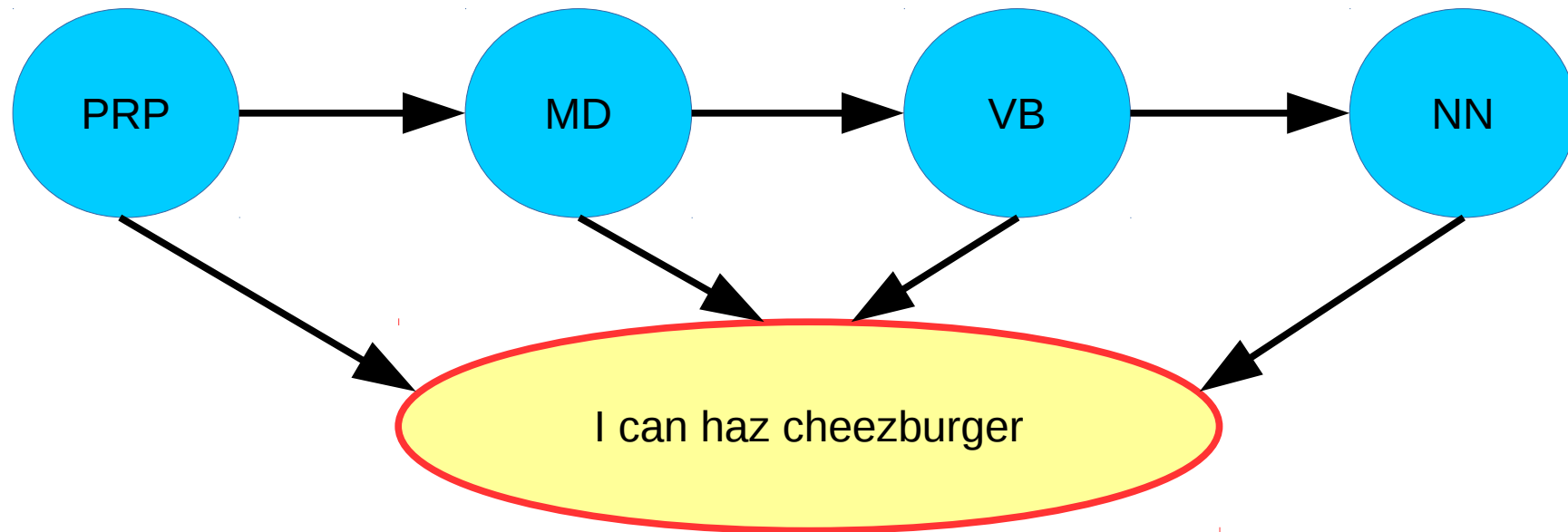
Statistical NLP

- Hidden Markov Models (HMMs)



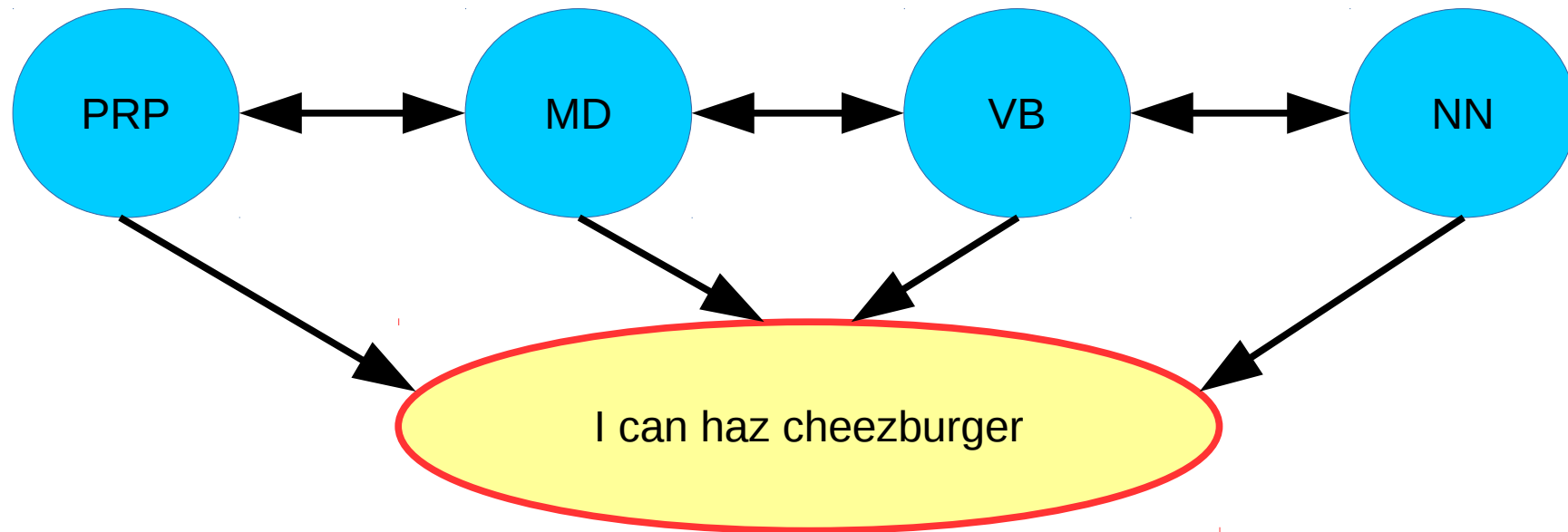
Statistical NLP

- Hidden Markov Models (HMMs); **MEMMs**



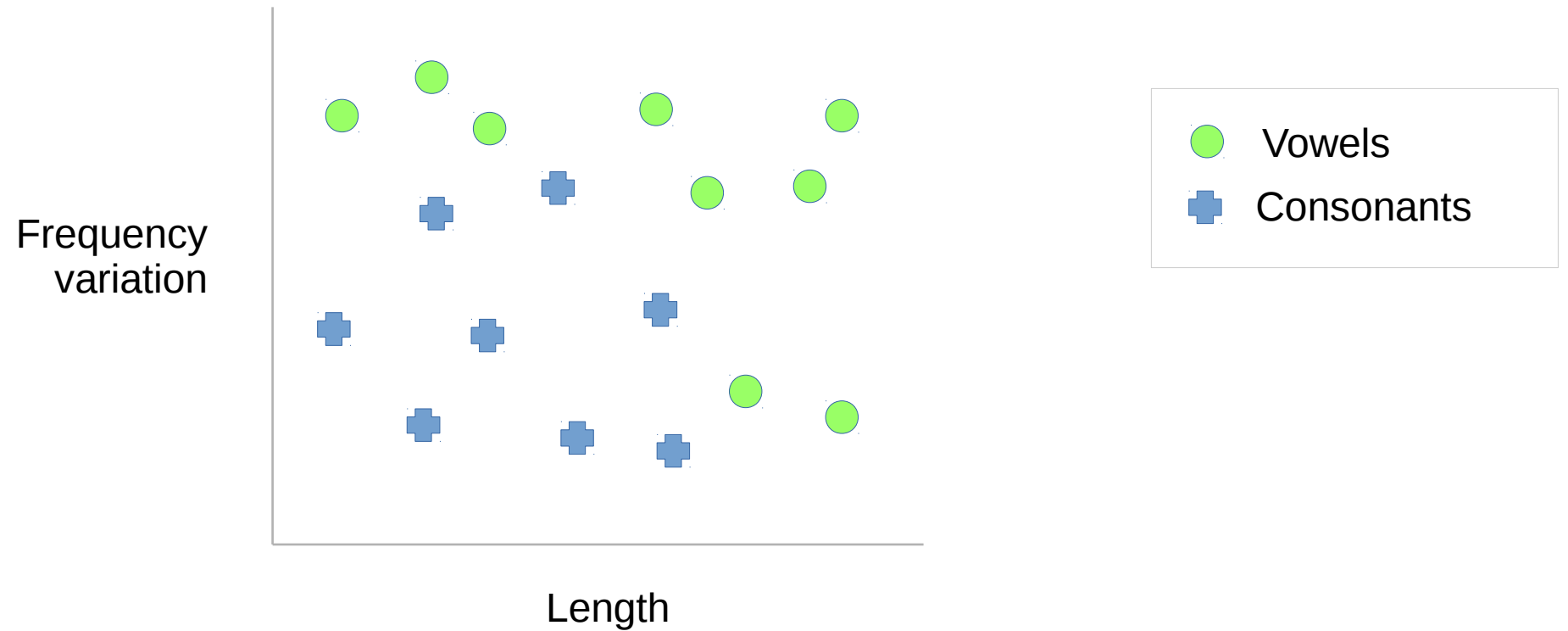
Statistical NLP

- Hidden Markov Models (HMMs); MEMMs, **CRFs**



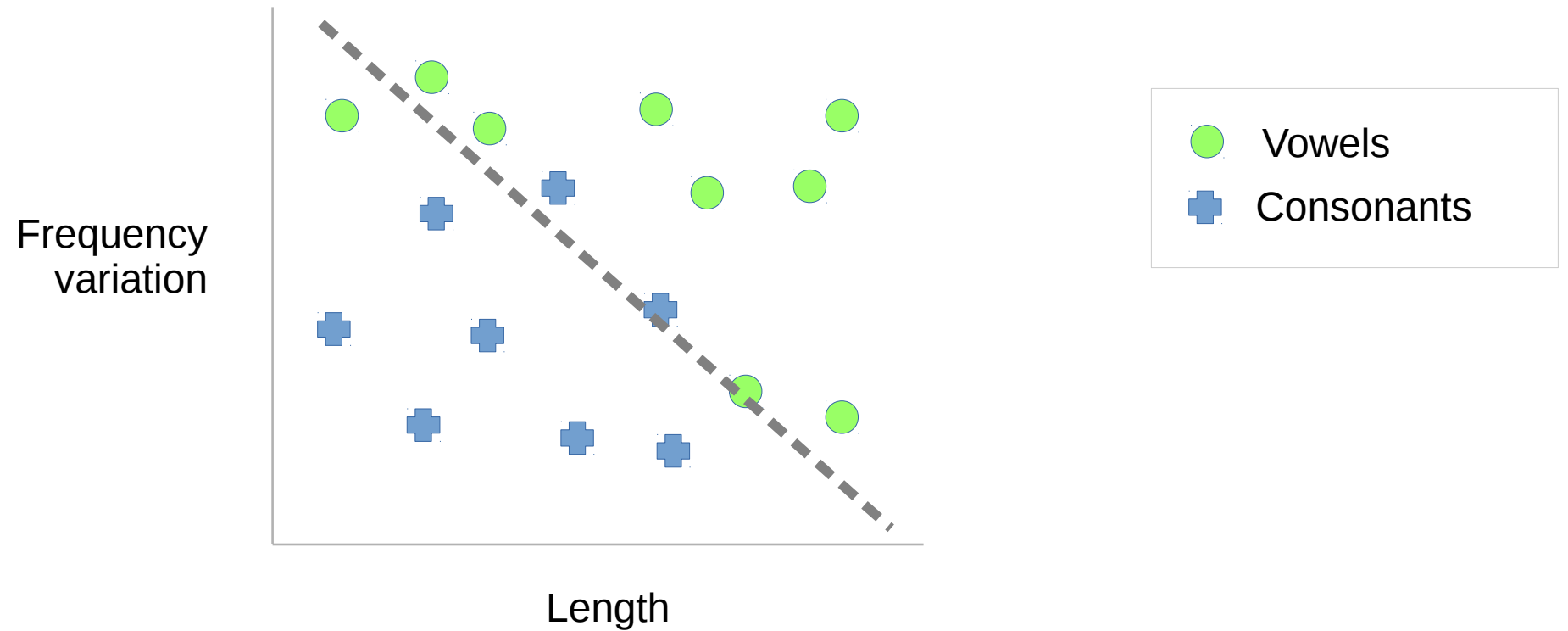
Statistical NLP

- Support Vector Machines (SVMs)



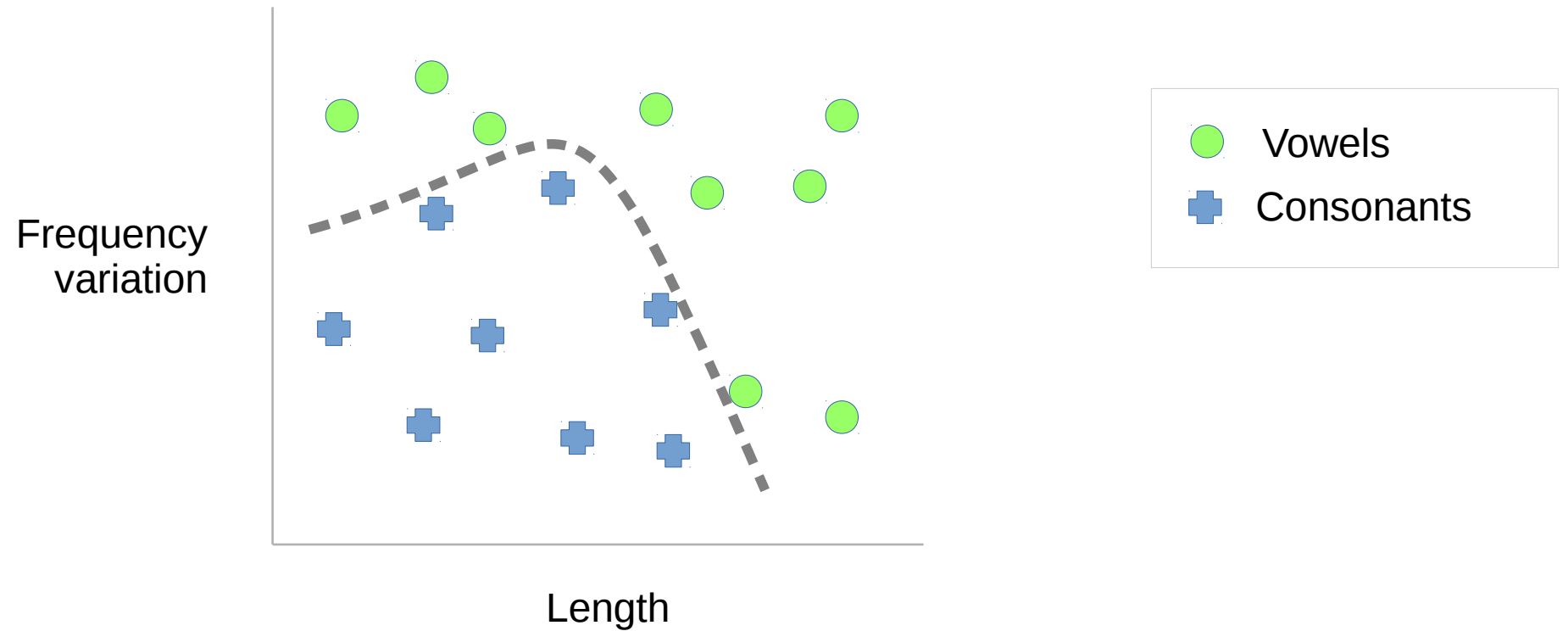
Statistical NLP

- Support Vector Machines (SVMs)



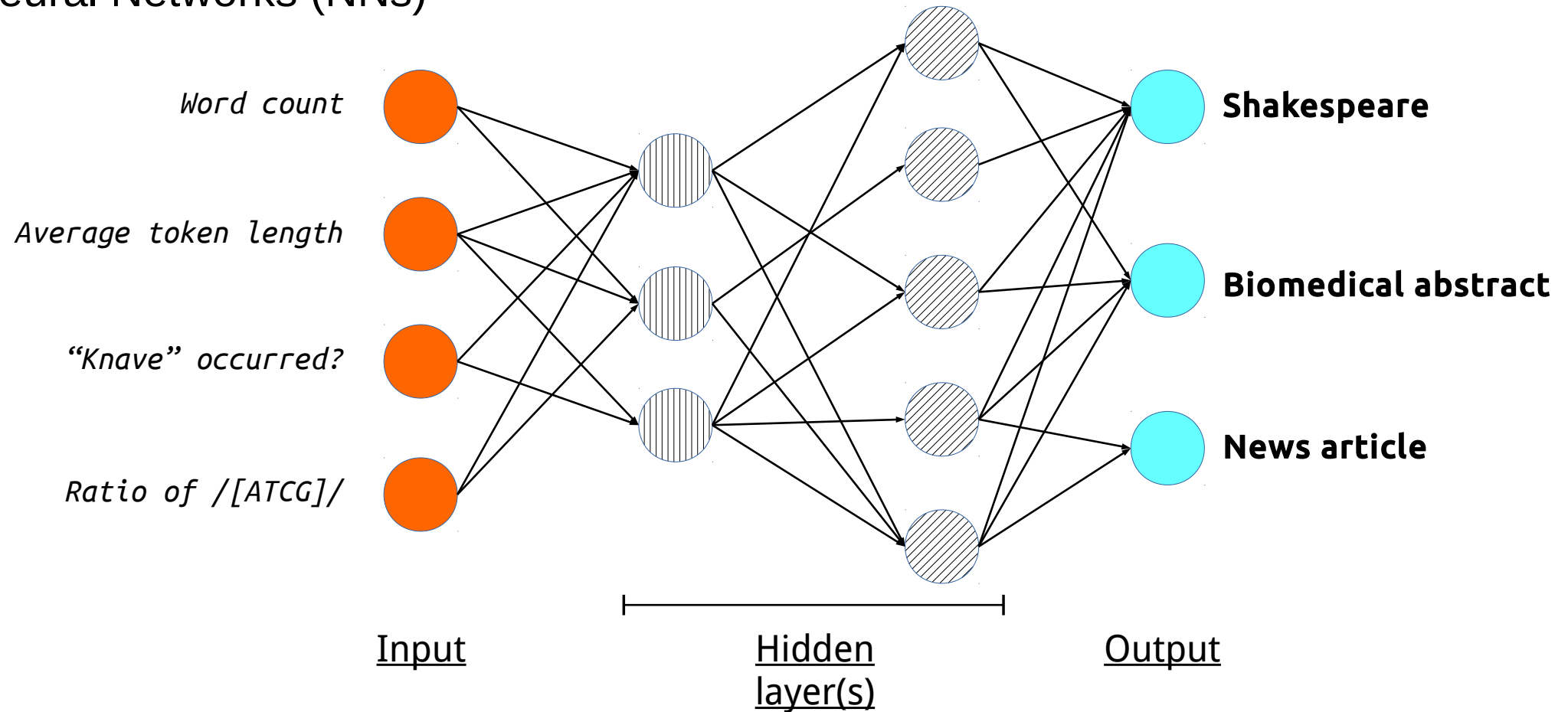
Statistical NLP

- Support Vector Machines (SVMs)



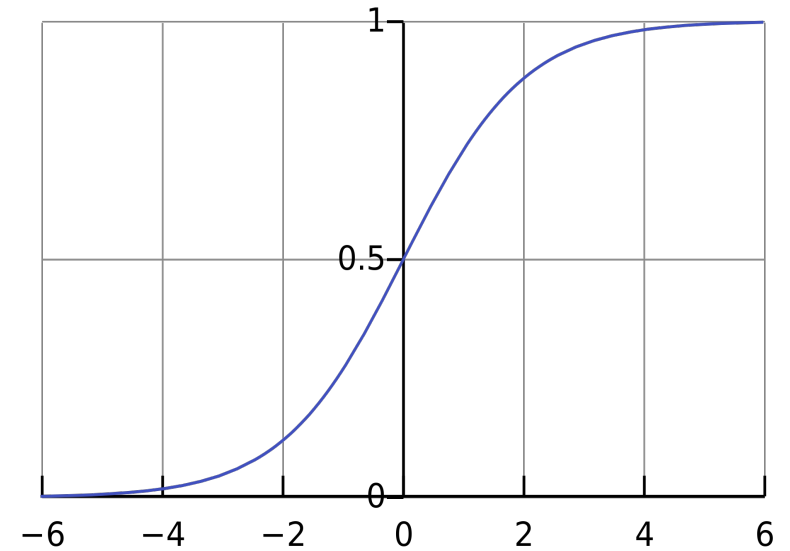
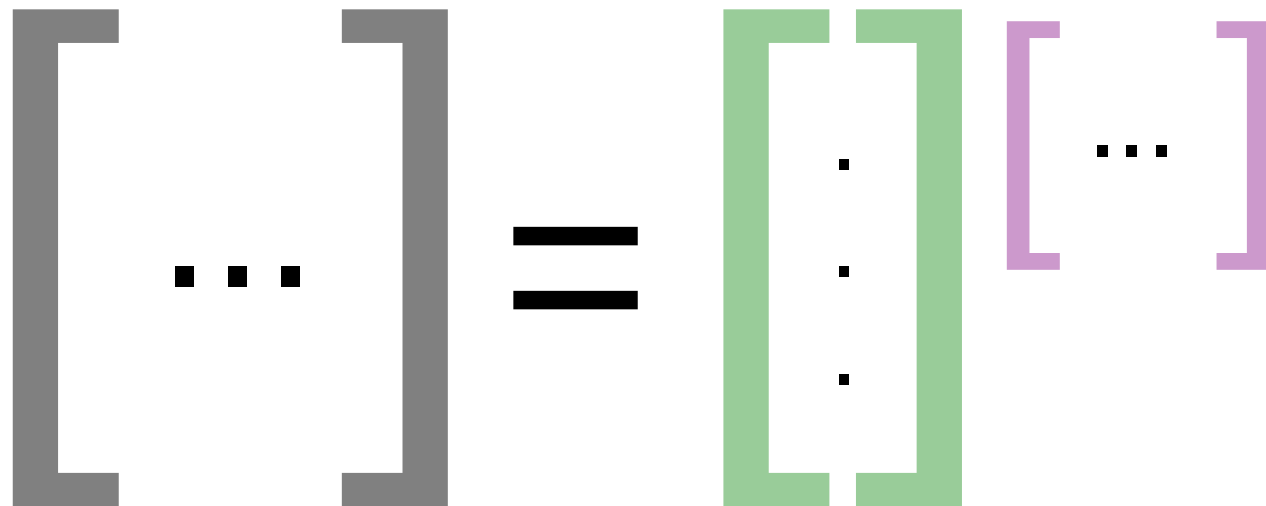
Statistical NLP

- Neural Networks (NNs)



Statistical NLP

- Other methods: matrix factorization, logistic regression, etc.



~~Rule-Based~~ NLP

~~Statistical~~ NLP

Lots of current work
uses both approaches
in **joint systems!**

These are models...

These are models...

...but models are only tools to solve
problems.

Kinds of Machine Learning

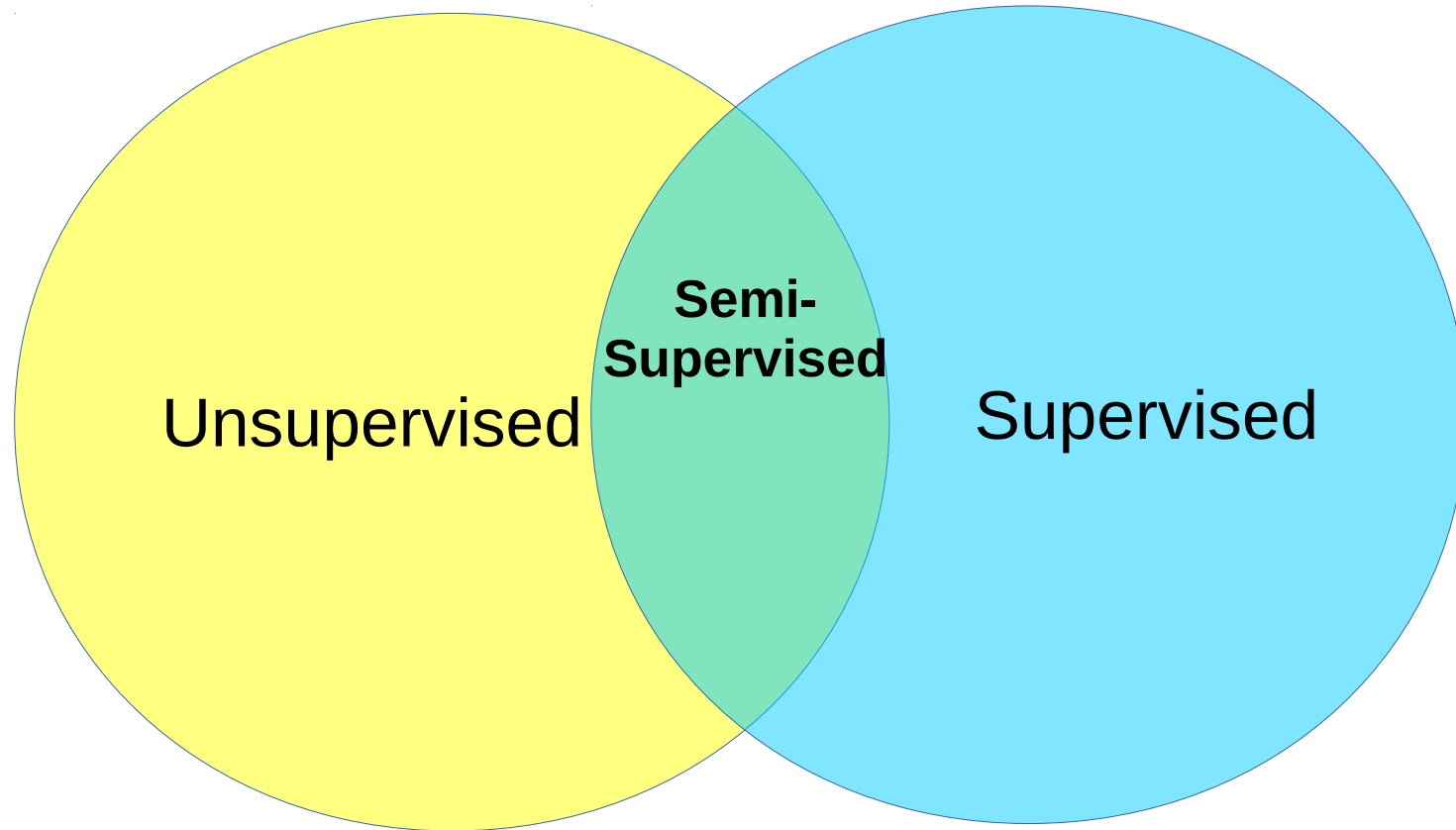


Unsupervised

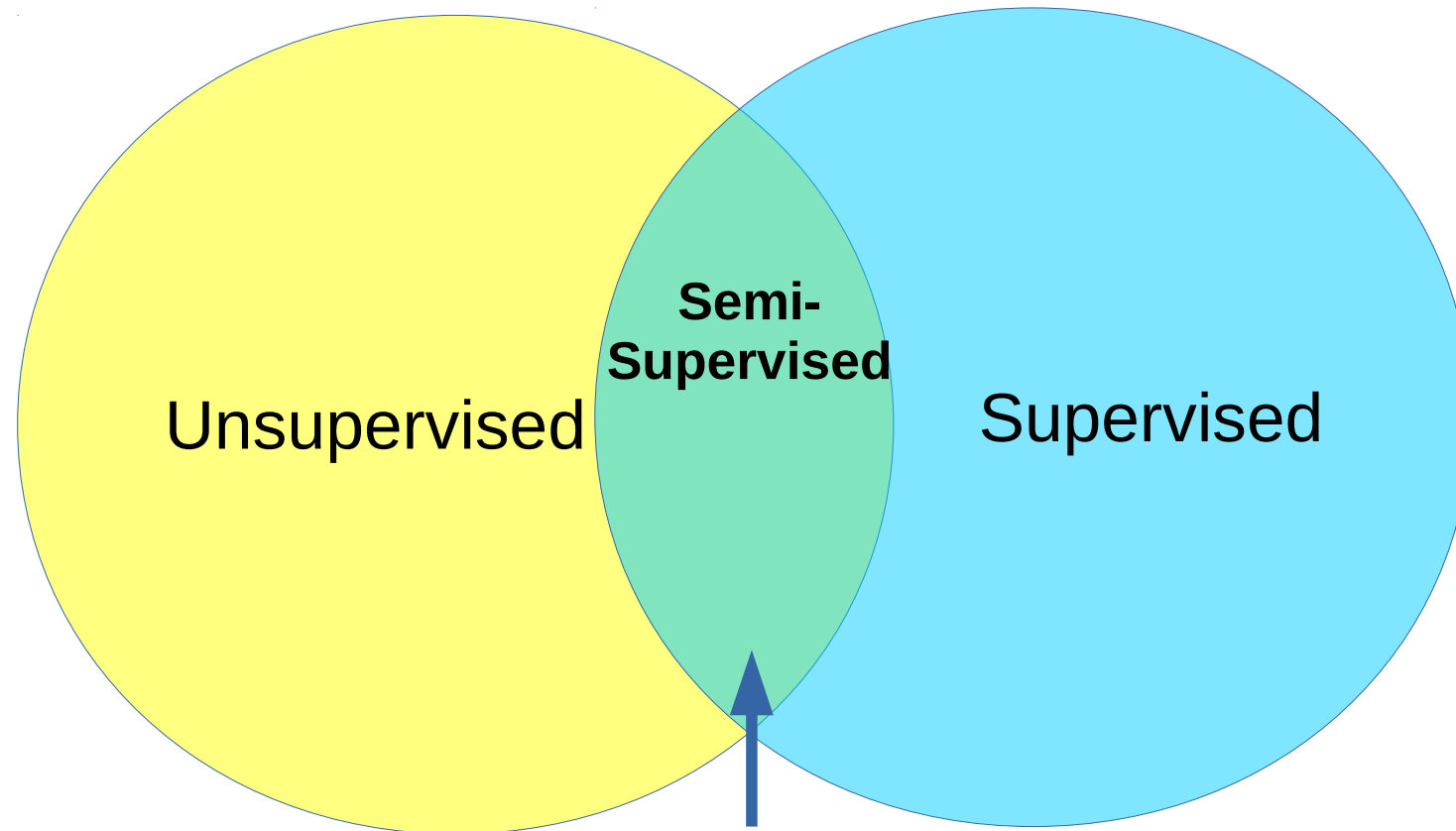


Supervised

Kinds of Machine Learning



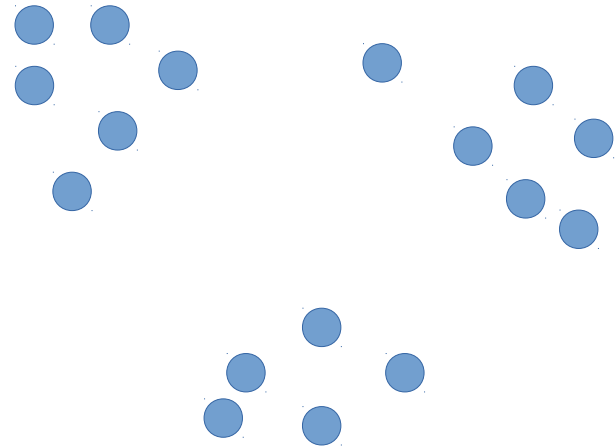
Kinds of Machine Learning



*aka Distantly-supervised,
weakly-supervised*

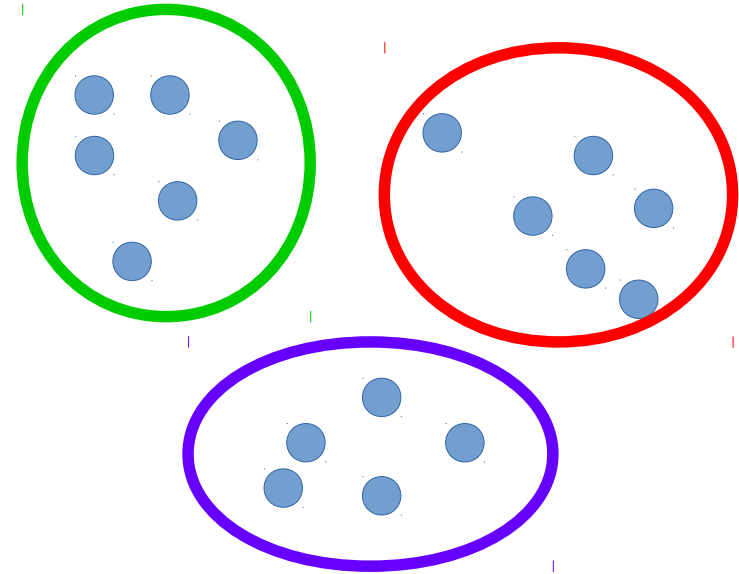
Unsupervised Learning

Goal: Discover hidden structure in data

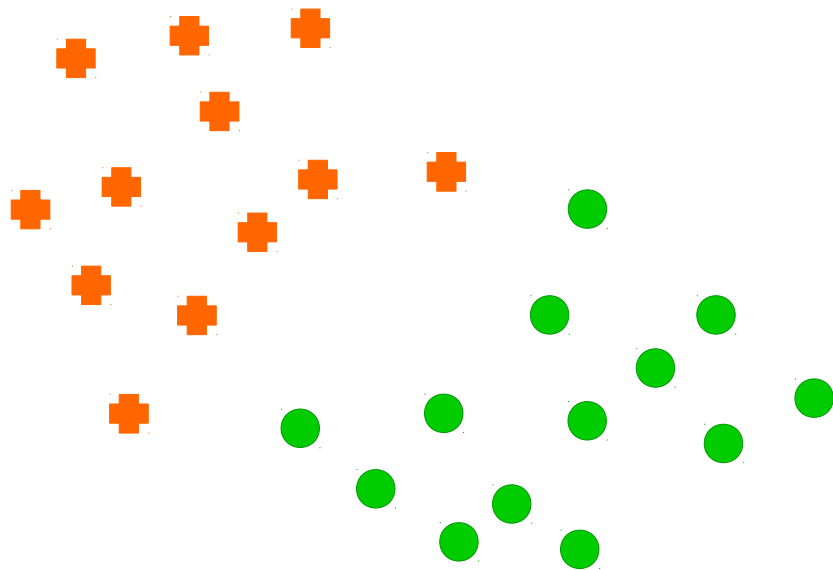


Unsupervised Learning

Goal: Discover hidden structure in data

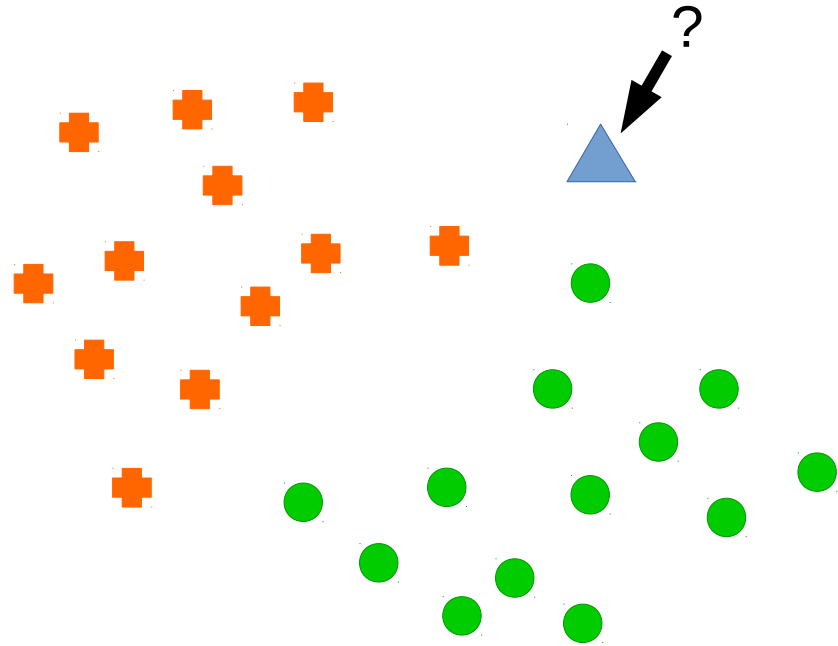


Supervised Learning



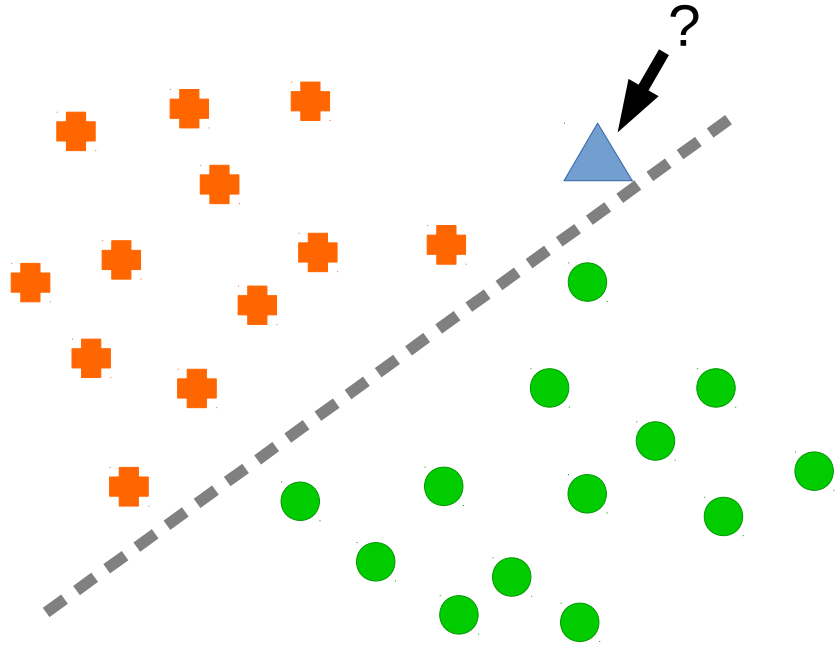
Goal: Use known information to categorize data

Supervised Learning



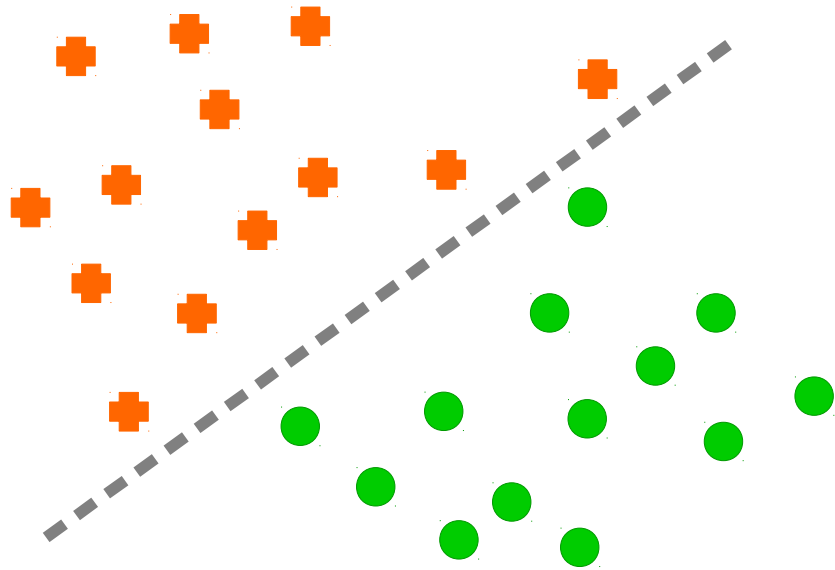
Goal: Use known information to categorize data

Supervised Learning



Goal: Use known information to categorize data

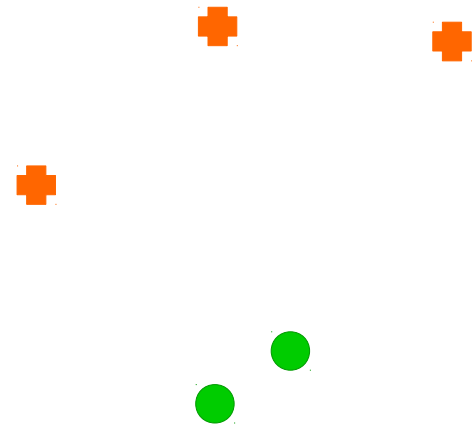
Supervised Learning



Goal: Use known information to categorize data

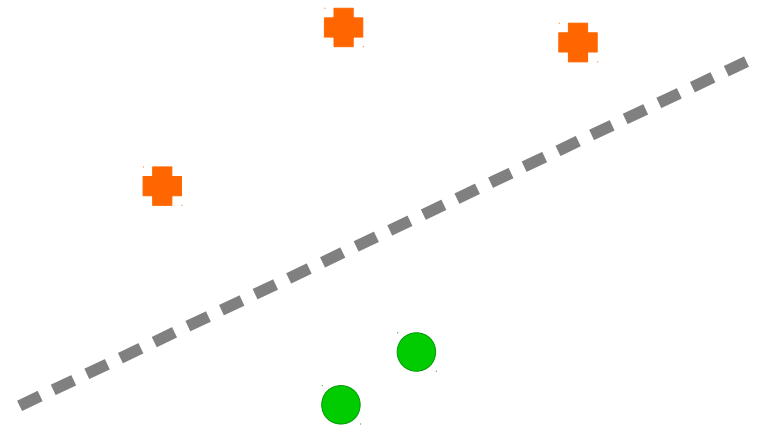
Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



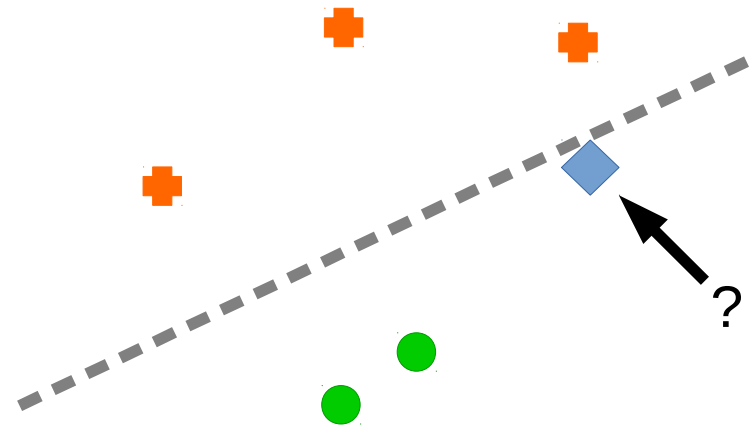
Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



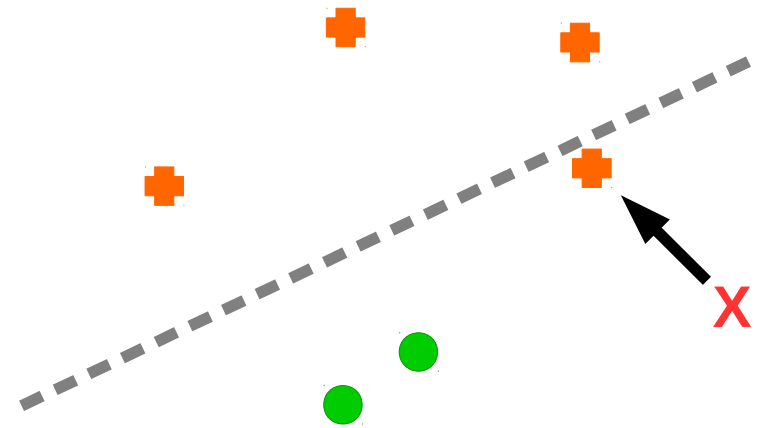
Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



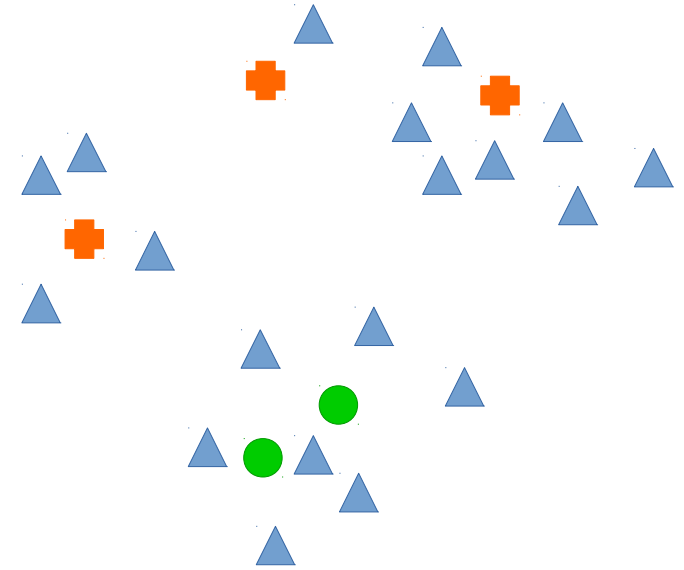
Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



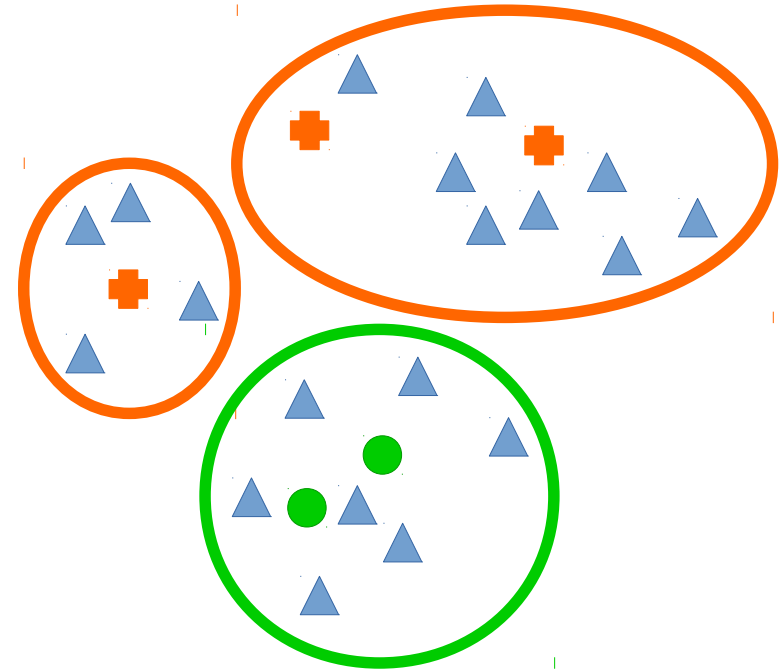
Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



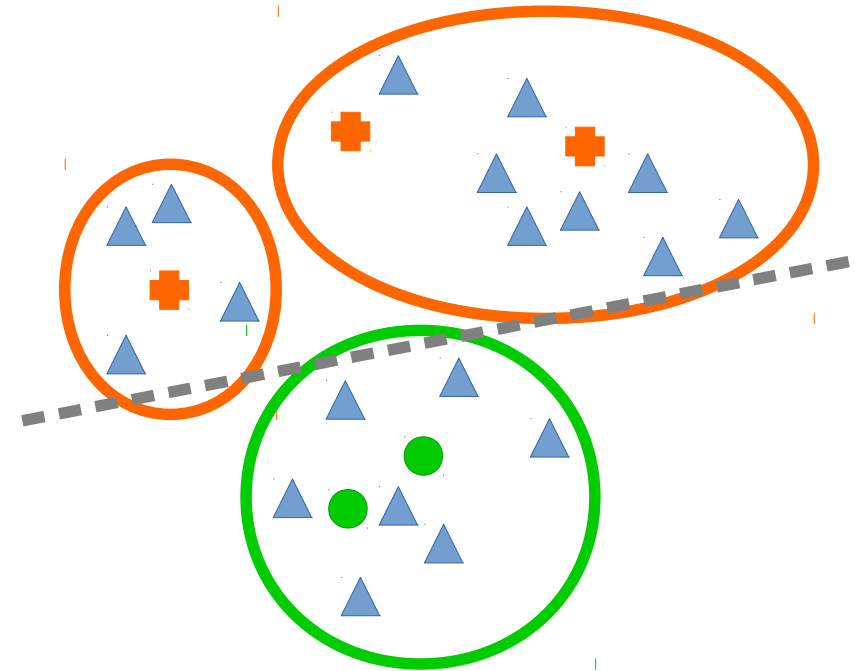
Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



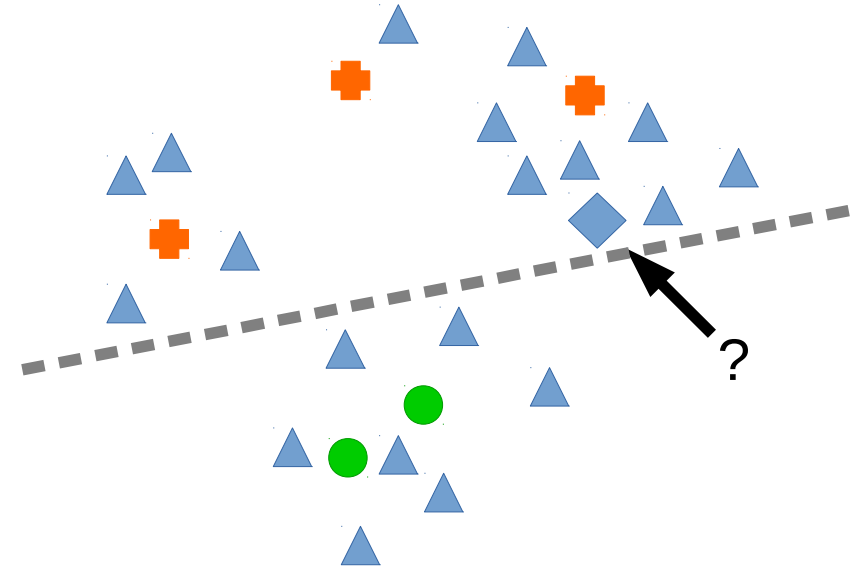
Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



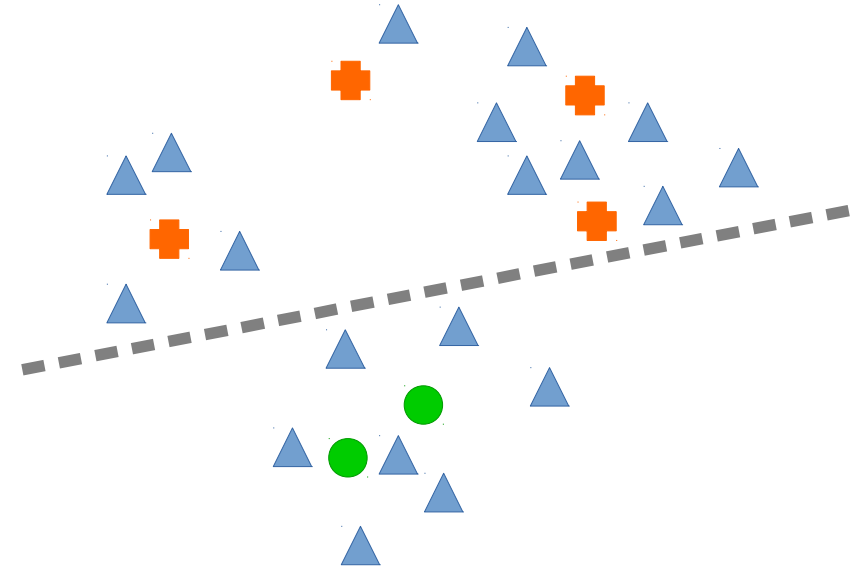
Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



Semi-Supervised Learning

Goal: Use some known information, along with hidden structure, to categorize data



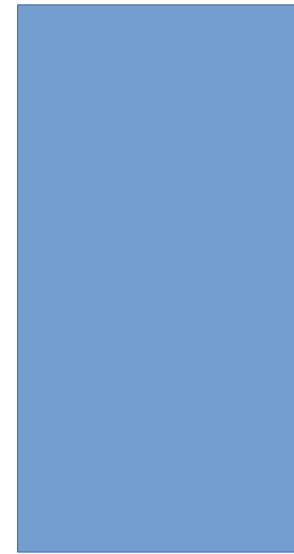
Human effort required

Human effort required

Unsupervised

Human effort required

Unsupervised



Supervised

Human effort required



Unsupervised

Semi-supervised

Supervised

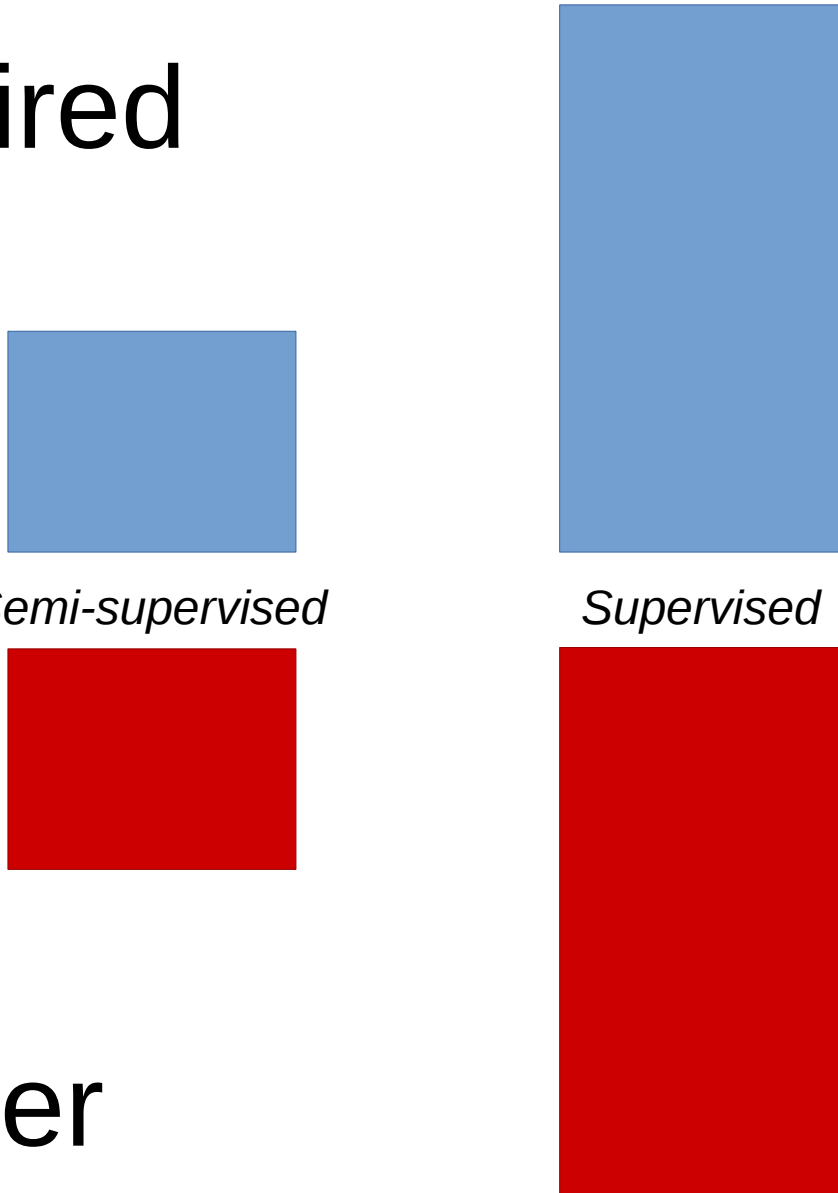
Human effort required

Unsupervised

Semi-supervised

Supervised

Classification power



At this point, you may be
asking yourself...

At this point, you may be asking yourself...

?

?

?

?

So what do you **do** with all this stuff?

?

?

?

Lots of things!

Machine Translation

Translate



English Spanish French English - detected



English French Russian

Translate

This translation |sucks

Этот перевод отстой



Wrong?

Etot perevod ostoy

Parsing / Tagging

Picard

ordered

tea.

Parsing / Tagging

Picard

ordered

tea.

Part of Speech

NNP

VBD

NN

Parsing / Tagging

Picard

ordered

tea.

Part of Speech

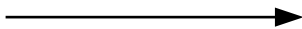
NNP

VBD

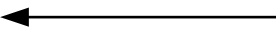
NN

Dependency

NSubj



Root



DObj

Information Extraction

*“Abraham Lincoln was born
February 12, 1809, in Hardin
County, Kentucky...”*

Information Extraction

*“**Abraham Lincoln** was born February 12, 1809, in Hardin County, Kentucky...”*



Birth Dates

| ID | Month | Day | Year |
|------------|----------|-----|------|
| Honest Abe | February | 12 | 1809 |

Birth Locations

| ID | County | State | Country |
|-------------|--------|----------|---------|
| Big Lincoln | Hardin | Kentucky | 'Murica |

Information Retrieval

Information Retrieval

*Web
Search*



Information Retrieval

*Web
Search*

Google

bing

Bioinformatics

ATTACCGCAGAT



1 | CATTACCGGAGATCCTA
2 | CCCATTACGGCCGCAGATAA
3 | ATTACCGAA

Information Retrieval

*Web
Search*

Google

bing

Bioinformatics

ATTACCGCAGAT



1 | CATTACCGGAGATCCTA
2 | CCCATTACGGCCGCAGATAA
3 | ATTACCGAA

*Question
Answering*

Who played Malcolm
Reynolds?

Nathan Fillion

Who played Real Madrid
last week?

Barcelona; final score 3-2

Etc., etc., etc.

Etc., etc., etc.

Automatic summarization

Etc., etc., etc.

Automatic summarization

Bacon ipsum dolor amet spare ribs leberkas filet mignon t-bone tenderloin ground round. Leberkas kevin meatball, short ribs rump andouille meatloaf pancetta shank bacon pork belly frankfurter picanha shankle sausage. Salami strip steak sirloin cow. Andouille ball tip meatloaf biltong bresaola. Cupim drumstick swine t-bone pork belly frankfurter jowl chuck leberkas cow short ribs ball tip.

Porchetta leberkas swine kevin ham capicola shankle strip steak hamburger salami filet mignon tri-tip bresaola picanha. Brisket tail swine biltong, capicola shankle sirloin. Jerky meatloaf ribeye, fatback turkey pork chop porchetta landjaeger ham salami meatball tongue pancetta kevin. Tri-tip swine filet mignon meatloaf bresaola porchetta pancetta salami frankfurter pork chop. Pork loin jerky pork chop, drumstick chuck flank ground round. Landjaeger hamburger pastrami salami.

Etc., etc., etc.

Automatic summarization

Bacon ipsum dolor amet spare ribs leberkas filet mignon t-bone tenderloin ground round. Leberkas kevin meatball, short ribs rump andouille meatloaf pancetta shank bacon pork belly frankfurter picanha shankle sausage. Salami strip steak sirloin cow. Andouille ball tip meatloaf biltong bresaola. Cupim drumstick swine t-bone pork belly frankfurter jowl chuck leberkas cow short ribs ball tip.

Porchetta leberkas swine kevin ham capicola shankle strip steak hamburger salami filet mignon tri-tip bresaola picanha. Brisket tail swine biltong, capicola shankle sirloin. Jerky meatloaf ribeye, fatback turkey pork chop porchetta landjaeger ham salami meatball tongue pancetta kevin. Tri-tip swine filet mignon meatloaf bresaola porchetta pancetta salami frankfurter pork chop. Pork loin jerky pork chop, drumstick chuck flank ground round. Landjaeger hamburger pastrami salami.



**Bacon bacon bacon
bacon pork!**

Etc., etc., etc.

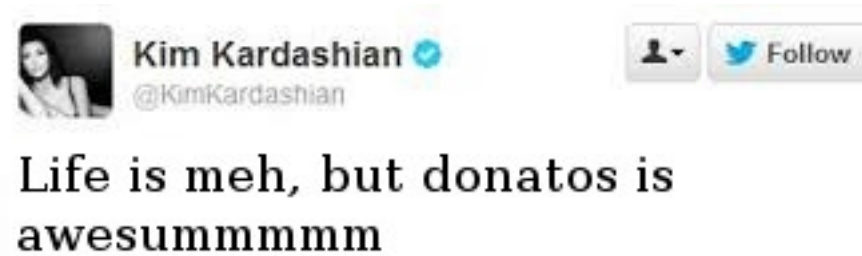
Automatic summarization

Sentiment analysis

Etc., etc., etc.

Automatic summarization

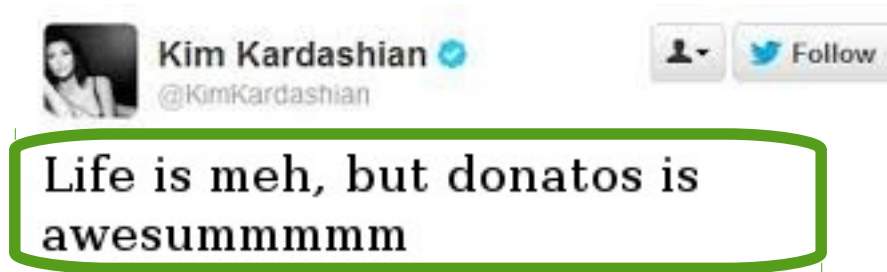
Sentiment analysis



Etc., etc., etc.

Automatic summarization

Sentiment analysis



Etc., etc., etc.

Automatic summarization

Sentiment analysis



Etc., etc., etc.

Automatic summarization

Sentiment analysis

Discourse analysis

Etc., etc., etc.

Automatic summarization

U: I want Chinese food.

S: Here are 473 Chinese places.

Sentiment analysis

U: How about cheap ones on the south side?

S: Here is 1 restaurant.

Discourse analysis

U: Eh, let's do Thai food instead.

S: I'm sorry, Dave, I can't let you do that.

Etc., etc., etc.

Automatic summarization

Sentiment analysis

Discourse analysis

U: I want Chinese food.

S: Here are 473 Chinese places.

U: How about cheap ones on the south side?

S: Here is 1 restaurant.

U: Eh, let's do Thai food instead.

S: I'm sorry, Dave, I can't let you do that.

User Goals

| <i>Turn</i> | <i>Type</i> | <i>Location</i> | <i>Cheap?</i> |
|-------------|-------------|-----------------|---------------|
| 1 | Chinese | ??? | ??? |
| 2 | Chinese | South | Yes |
| 3 | Thai | South | Yes |

Etc., etc., etc.

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Etc., etc., etc.

Phonemes

U|n|b|r|ea|k|a|b|le

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Etc., etc., etc.

Phonemes

U|n|b|r|ea|k|a|b|le

Automatic summarization

Morphemes

Sentiment analysis

Un|break|able

Discourse analysis

Segmentation

Etc., etc., etc.

Phonemes

U|n|b|r|ea|k|a|b|le

Automatic summarization

Morphemes

Un|break|able

Sentiment analysis

Discourse analysis

Words

maytheforcebewithyou

Segmentation



May the force be with you

Etc., etc., etc.

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Phonemes

U|n|b|r|ea|k|a|b|le

Morphemes

Un|break|able

Words

maytheforcebewithyou



May the force be with you

Sentences

[I spoke to Mr. Spock.]
[His response was
illogical.]

Etc., etc., etc.

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Phonemes

U|n|b|r|ea|k|a|b|le

Morphemes

Un|break|able

Words

maytheforcebewithyou



May the force be with you

Sentences

[I spoke to Mr. Spock.]
[His response was
illogical.]

Topics

...who I met at a
Trek convention.

As for Star Wars...

Etc., etc., etc.

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Disambiguation and reference

Word sense disambiguation

Etc., etc., etc.

After I put him in [check]¹, he wrote
me a [check]².

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Disambiguation and reference

Etc., etc., etc.

Word sense disambiguation

After I put him in [check]¹, he wrote
me a [check]².

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Disambiguation and reference

Coreference resolution

I spoke to [the customer]₁, then told [my
boss]₂ that [she]₂ should fire [her]₁.

Etc., etc., etc.

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Disambiguation and reference

Word sense disambiguation

After I put him in [check]¹, he wrote
me a [check]².

Coreference resolution

I spoke to [the customer]₁, then told [my
boss]₂ that [she]₂ should fire [her]₁.

Named entity recognition

[Bugs Bunny]_{Person} bought 50% of
[Acme Corp.]_{Company} in [2004]_{Year}.

Etc., etc., etc.

Automatic summarization

Sentiment analysis

Discourse analysis

Segmentation

Disambiguation and reference

And many more!

How can I get in on this?

NLP Toolkits

| <i>Toolkit</i> | <i>Language</i> | <i>Website</i> |
|---|-----------------|---|
| Apache OpenNLP General-purpose NLP toolkit; tends to use older models, but under Apache license. | Java | https://opennlp.apache.org |
| Natural Language Toolkit (NLTK) Standard NLP option for Python; easy to pick up and play with, and includes several common corpora. | Python | http://www.nltk.org/ |
| Mallet More technical toolkit, focused on current, high-complexity models. | Java | http://mallet.cs.umass.edu/ |
| LingPipe Another general-purpose NLP toolkit; offers industry licensing option. | Java | http://alias-i.com/lingpipe/ |
| Stanford CoreNLP Standard tools in academia, tends towards cutting edge models. Low ease-of-use, and academic licensing restrictions. | Java | http://nlp.stanford.edu/software/corenlp.shtml |
| Alchemy API Fanciest industry option (owned by IBM). Offers NLP, vision, other ML resources. | Cloud API | http://www.alchemyapi.com/ |

Other Resources



Speech Recognition Toolkit - <http://kaldi-asr.org/>



<http://www.signalprocessingsociety.org/>



Association for Computational Linguistics

<http://aclweb.org/>

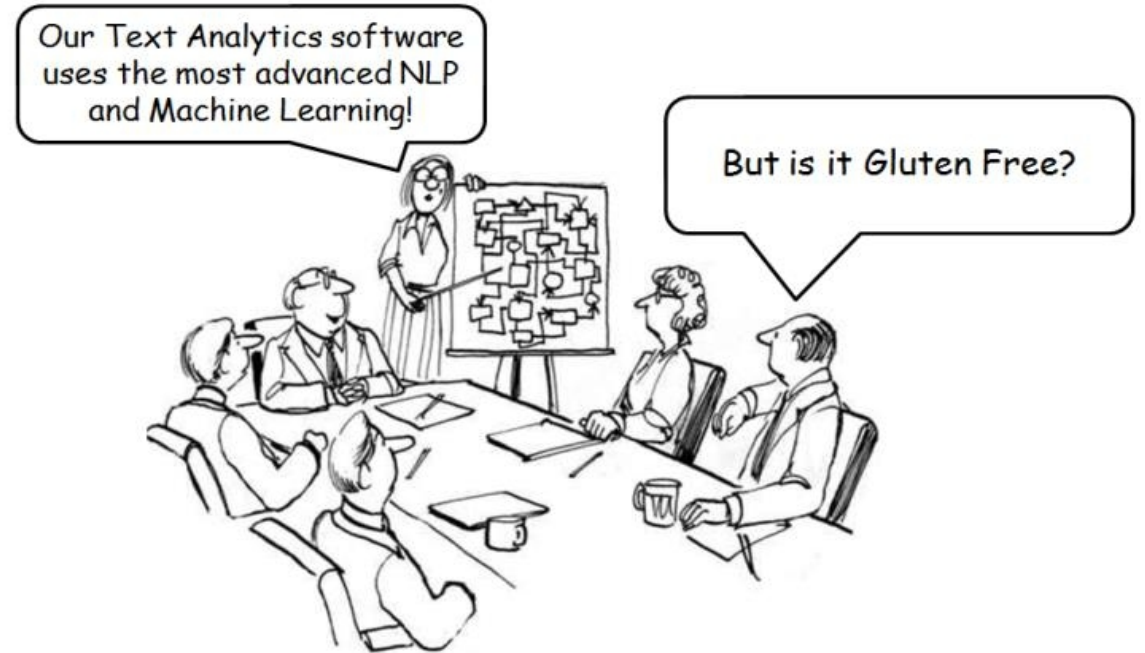
Questions?

My contact info:

Denis Griffis

griffis.30@osu.edu

<http://web.cse.ohio-state.edu/slate/>



MBA Rule #1:
Always Counter Buzz Words with Buzz Words

@TomHCAnderson
tomhcanderson.com