

# Natural Language Processing

Nuri Cingillioglu

<https://www.doc.ic.ac.uk/~nuric/>

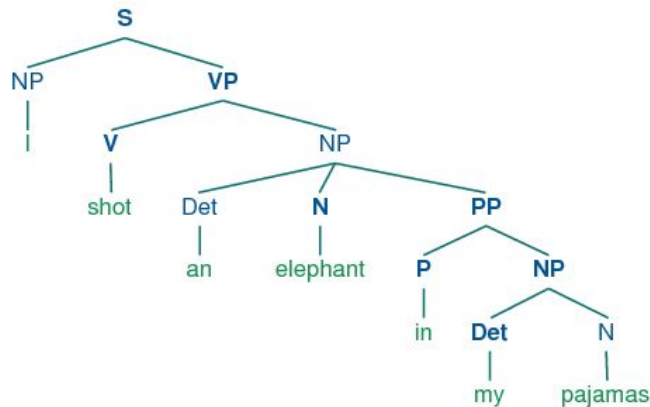
Many thanks to Lucia Specia

# Dependency parsing

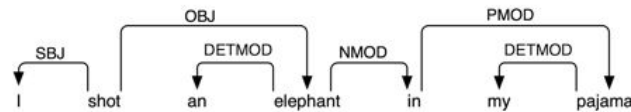
# Definition

- Given a sequence of words (usually a sentence), generate its **syntactic structure**

## Constituency Parsing



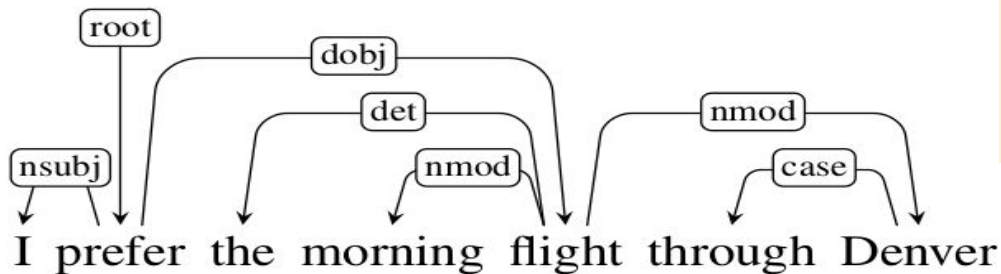
## Dependency Parsing



# Dependency parsing

Figure from Jurafsky, D and Martin, J, “Speech and Language Processing,” 2018, ch 15

- Connect words in sentence to indicate dependencies between them - much older linguistic theory
- Build around notion of having **heads** and **dependents**
- Arrows can be annotated by different types of dependencies
  - **Head** (governor), also called **argument**: origin
  - **Dependent**, also called **modifier**: destiny



From head  
to dependent

# Dependency parsing

- There are versions without **typed dependencies**, just arcs

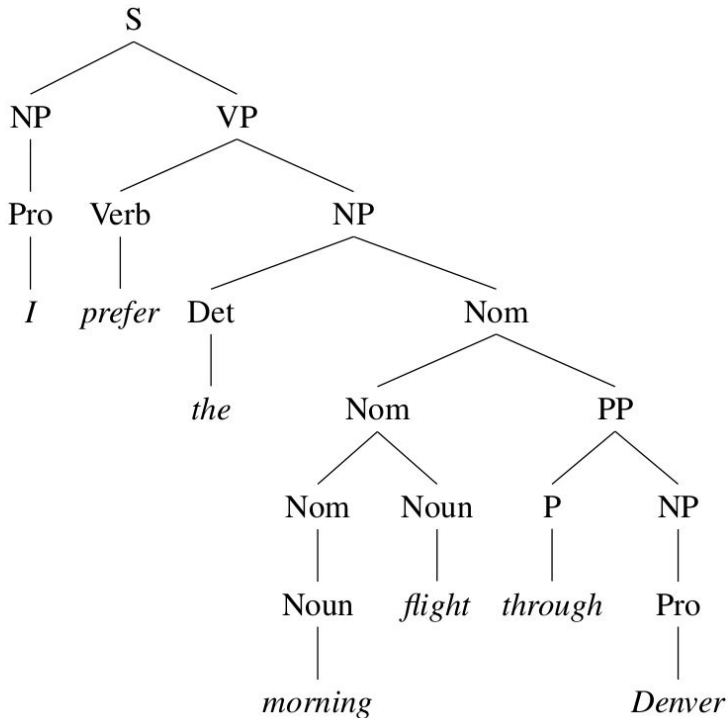
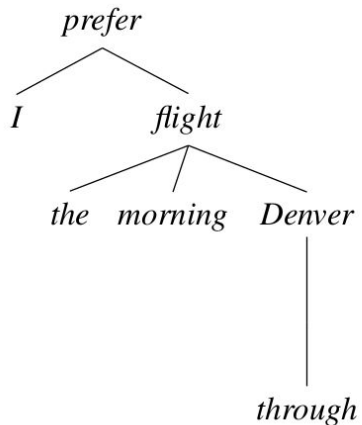


- Untyped variant is simpler to build but less informative
- Can also represent as a tree

# Dependency parsing

Figure from Jurafsky, D and Martin, J, “Speech and Language Processing,” 2018, ch 15

- Comparison to constituency parsing



# Dependency parsing

- Main **differences** to constituency parsing
  - No nodes corresponding to phrasal constituents or lexical categories
  - The internal structure consists only of directed relations between **pairs** of lexical items
  - These relationships allow directly encoding important information, e.g.:
    - The arguments of the verb *prefer* are directly linked to it in the dependency structure

# Dependency parsing

- Main **advantages** in dependency parsing
  - Ability to deal with languages that are morphologically rich and have a relatively free word order. E.g. Czech location adverbs may occur before or after object:  
I caught the fish **here** vs I caught **here** the fish
  - Would have to represent two rules for each possible place of the adverb for constituency
  - Dependency approach: only one link; abstracts away from word order



# Dependency parsing

- Main **advantages** in dependency parsing
  - Head-dependent relations provide approximation to semantic relationship between predicates and arguments
  - Can be directly used to solve problems such as co-reference resolution, question answering, etc.

# Dependency relations (e.g. from universal dep.)

Figure from Jurafsky, D and Martin, J, “Speech and Language Processing,” 2018, ch 15

<b>Clausal Argument Relations</b>	<b>Description</b>
NSUBJ	Nominal subject
DOBJ	Direct object
IOBJ	Indirect object
CCOMP	Clausal complement
XCOMP	Open clausal complement
<b>Nominal Modifier Relations</b>	<b>Description</b>
NMOD	Nominal modifier
AMOD	Adjectival modifier
NUMMOD	Numeric modifier
APPOS	Appositional modifier
DET	Determiner
CASE	Prepositions, postpositions and other case markers
<b>Other Notable Relations</b>	<b>Description</b>
CONJ	Conjunct
CC	Coordinating conjunction

# Dependency relations (e.g. from universal dep.)

Figure from Jurafsky, D and Martin, J, "Speech and Language Processing," 2018, ch 15

Relation	Examples with <i>head</i> and <b>dependent</b>
NSUBJ	<b>United</b> <i>canceled</i> the flight.
DOBJ	United <i>diverted</i> the <b>flight</b> to Reno. We <i>booked</i> her the first <b>flight</b> to Miami.
IOBJ	We <i>booked</i> <b>her</b> the flight to Miami.
NMOD	We took the <b>morning</b> <i>flight</i> .
AMOD	Book the <b>cheapest</b> <i>flight</i> .
NUMMOD	Before the storm JetBlue canceled <b>1000</b> <i>flights</i> .
APPOS	<i>United</i> , a <b>unit</b> of UAL, matched the fares.
DET	<b>The</b> <i>flight</i> was canceled. <b>Which</b> <i>flight</i> was delayed?
CONJ	We <i>flew</i> to Denver and <b>drove</b> to Steamboat.
CC	We flew to Denver <b>and</b> <i>drove</i> to Steamboat.
CASE	Book the flight <b>through</b> <i>Houston</i> .

# Dependency formalisms – general case

- A dependency structure is a **directed graph**  $G = (V, A)$ 
  - $V$  = set of vertices (words, punctuation, ROOT)
  - $A$  = set of ordered pairs of vertices (i.e. arcs)
- A dependency tree (directed graph):
  - Has a single ROOT node that has no incoming arcs
  - Each vertex has exactly one incoming arc (except for ROOT)
  - There's a unique path from ROOT to each vertex
  - There are no cycles  $A \rightarrow B, B \rightarrow A$

# Dependency formalisms – general case

- **This ensures the following properties:**
  - Dependency structure becomes a tree
  - Each word has a single head
  - The dependency tree is connected
  - There is a single ROOT from which a unique directed path follows to each word in sentence

# Dependency parsing –sources of info

- Distance between head and dependent
  - Mostly nearby words
- Intervening material
  - Dependencies don't cross over verbs or punctuation
- Valency of verbs
  - For a typical word, what kind of dependency it generally takes? E.g. A noun takes dependencies on the left (DET, JJ) but not on the right

# Dependency parsing – two approaches

- Dynamic programming (cubic time, not very accurate)
- Shift-reduce (transition-based)
  - Predict from left-to-right
  - Fast (linear), but slightly less accurate
  - **MaltParser**
- Spanning tree (graph-based, constraint satisfaction)
  - Calculate full tree at once
  - Slightly more accurate, slower
  - **MSTParser**

# Dependency parsing – transition-based

- Deterministic parsing, shift-reduce (MALT parser)
  - Greedy choice of attachment for each word in order, guided by ML classifier
  - Works very well in practice
  - **Linear time parsing!**



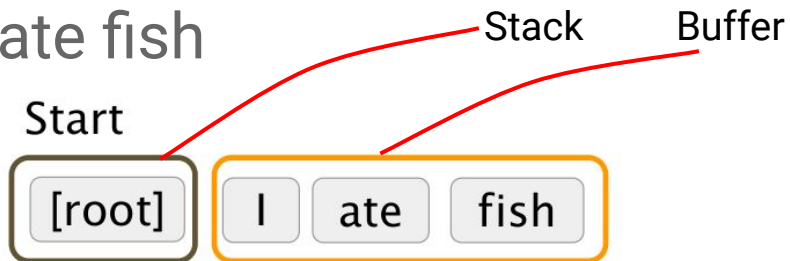
# Dependency parsing –MALT parser

- Reads sentence word by word, left to right
- Greedy decision as to how to attach each word as it is read
- Sequence of actions bottom-up
- Formally, 3 data structures:
  - $\sigma$  = **Stack**, which starts with ROOT
  - $\beta$  = **Buffer**, which starts with all words in sentence
  - $A$  = **Set of arcs**, which starts empty
- Set of actions
  - **Shift** / **left arc** / **right arc**
  - Optionally, set of **dep. labels** for left and right arc actions (~40)

# Dependency parsing –MALT parser

Example by Chris Manning

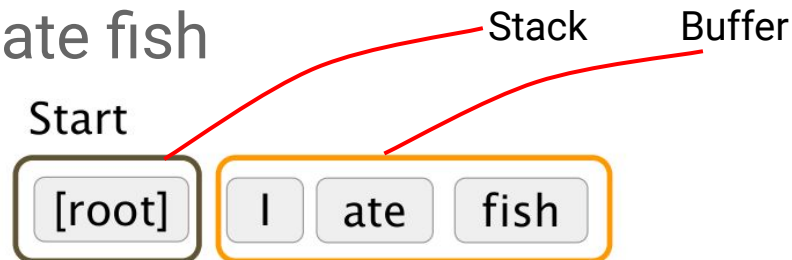
I ate fish



# Dependency parsing –MALT parser

Example by Chris Manning

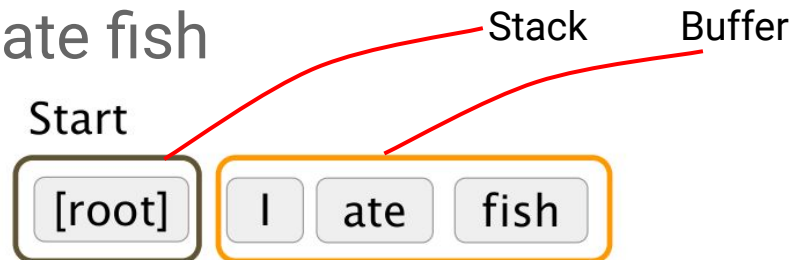
I ate fish



# Dependency parsing –MALT parser

Example by Chris Manning

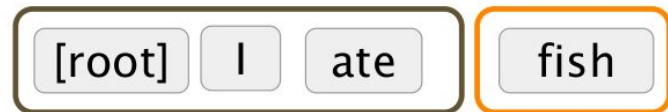
I ate fish



Shift



Shift



# Dependency parsing –MALT parser

Example by Chris Manning

Left Arc

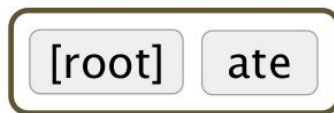


A +=  
nsubj(ate → I)

# Dependency parsing –MALT parser

Example by Chris Manning

Left Arc



A +=  
nsubj(ate → I)

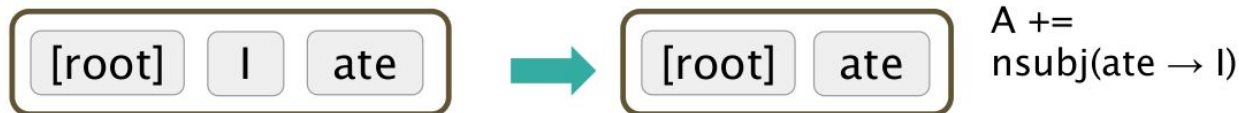
Shift



# Dependency parsing –MALT parser

Example by Chris Manning

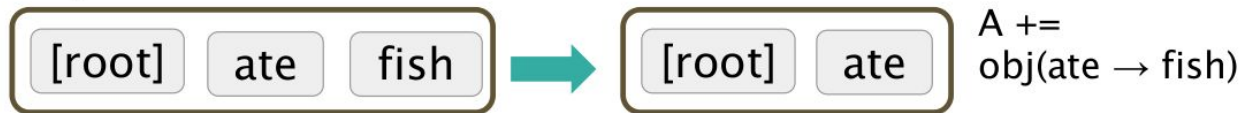
Left Arc



Shift



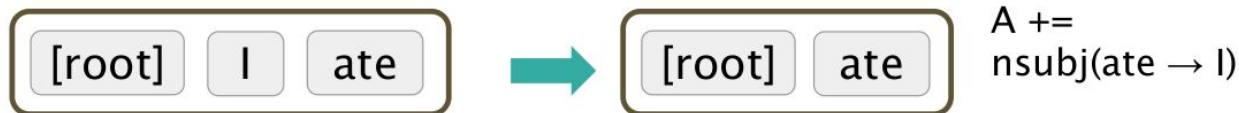
Right Arc



# Dependency parsing –MALT parser

Example by Chris Manning

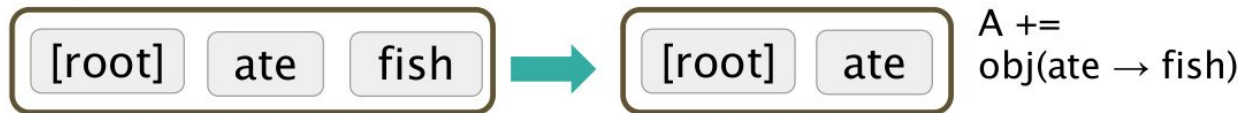
Left Arc



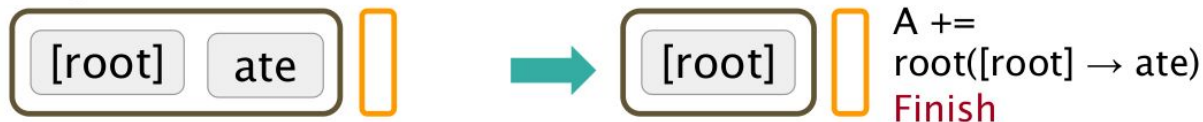
Shift



Right Arc



Right Arc





# Dependency parsing –MALT parser

Pseudocode by Chris Manning

**Start:**  $\sigma = [\text{ROOT}]$ ,  $\beta = w_1, \dots, w_n$ ,  $A = \emptyset$

1. Shift  $\sigma, w_i | \beta, A \Rightarrow \sigma | w_i, \beta, A$

2. Left-Arc<sub>r</sub>  $\sigma | w_i | w_j, \beta, A \Rightarrow \sigma | w_j, \beta, A \cup \{r(w_j, w_i)\}$

3. Right-Arc<sub>r</sub>  $\sigma | w_i | w_j, \beta, A \Rightarrow \sigma | w_i, \beta, A \cup \{r(w_i, w_j)\}$

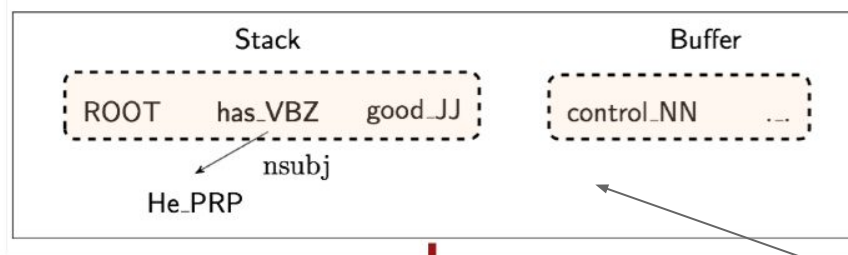
**Finish:**  $\sigma = [w]$ ,  $\beta = \emptyset$

# Dependency parsing –MALT parser

- How do we make the shift/reduce left/right decisions?
  - ML classifier
  - Each action is predicted by a **discriminative classifier** over each move
  - 3 classes for untyped dependencies: shift, left or right
  - $2 \times \text{categories} + 1$  for typed dependencies
  - **Features**: top of stack word, its POS; first in buffer word, its POS; etc.
  - No beam-search in original version

# Dependency parsing –MALT parser

Example by Chris Manning



encode

binary, sparse  
dim =  $10^6 \sim 10^7$

0 0 0 1 0 0 1 0 ... 0 0 1 0

Feature templates: usually a combination of 1 ~ 3 elements from the configuration.

Indicator features

$s1.w = \text{good} \wedge s1.t = \text{JJ}$   
 $s2.w = \text{has} \wedge s2.t = \text{VBZ} \wedge s1.w = \text{good}$   
 $lc(s_2).t = \text{PRP} \wedge s2.t = \text{VBZ} \wedge s1.t = \text{JJ}$   
 $lc(s_2).w = \text{He} \wedge lc(s_2).l = \text{nsubj} \wedge s2.w = \text{has}$

top of stack  
word, its POS;  
first in buffer  
word, its POS;  
etc.

# Dependency parsing – neural parser

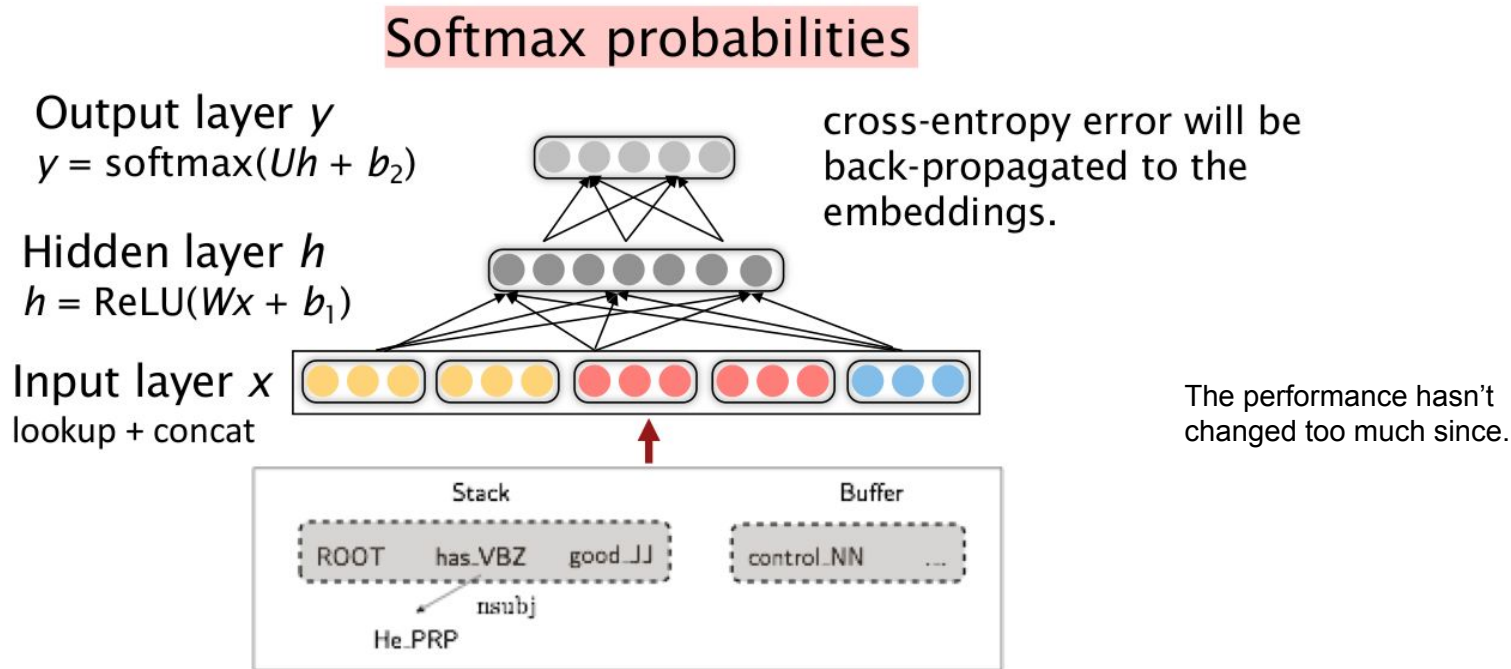
- Follow-up work (Chen and Manning, 2014)
  - Replace binary features by **embeddings** (words & POS)
  - Concatenate these embeddings
  - Train a FNN as classifier with cross-entropy loss
  - Superior performance!
    - UAS = untyped
    - LAS = typed dependencies

Parser	UAS	LAS	sent. / s
MaltParser	89.8	87.2	469
MSTParser	91.4	88.1	10
TurboParser	<b>92.3</b>	89.6	8
C & M 2014	92.0	<b>89.7</b>	<b>654</b>

# Dependency parsing – neural parser

Example by Chris Manning

- Follow-up work (Chen and Manning, 2014)

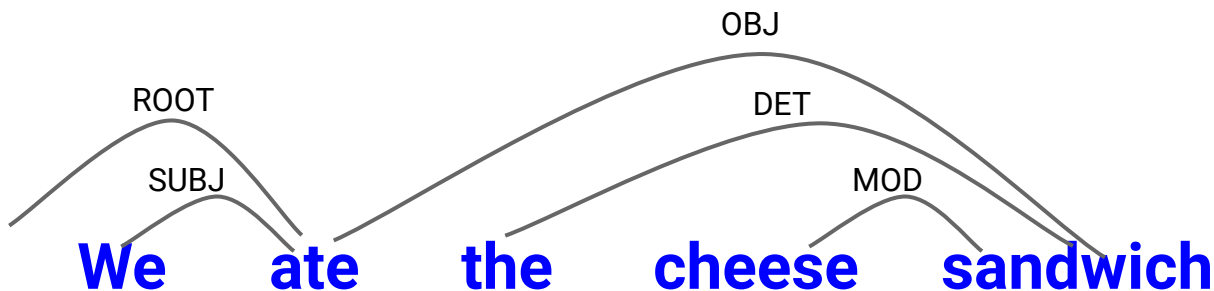


# Dependency parsing – treebanks

- Treebanks exist in the same way as for constituency parsing for training these classifiers
  - From converting constituency treebanks
  - Annotated from scratch, esp. For morphologically rich languages such as Czech, Hindi and Finnish

# Dependency parsing – evaluation

- Accuracy, precision or recall:
  - Count identical dependencies: span (non-typed parsing) or span and type of dependency (typed parsing). E.g.:



# Dependency parsing – evaluation

- Accuracy, precision or recall:
  - Count identical dependencies: span (non-typed parsing) or span and type of dependency (typed parsing). E.g.:

## Reference

(1,2) We	SUBJ
(2,0) eat	ROOT
(3,5) the	DET
(4,5) cheese	MOD
(5,2) sandwich	OBJ

## Hypothesis

(1,2) We	SUBJ
(2,0) eat	ROOT
(3,4) the	DET
(4,2) cheese	OBJ
(5,2) sandwich	PRED



# Dependency parsing – evaluation

- Accuracy, precision or recall:
  - Count identical dependencies: span (non-typed parsing) or span and type of dependency (typed parsing). E.g.:

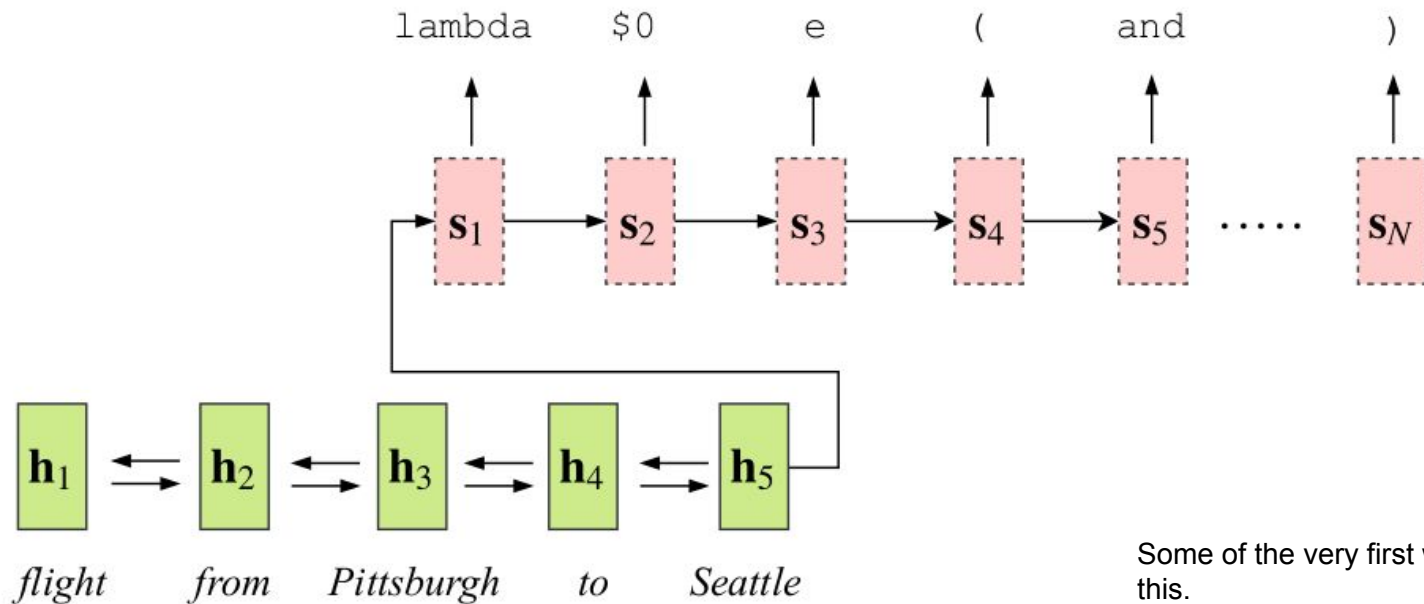
Reference		Hypothesis	
(1,2) We	SUBJ	(1,2) We	SUBJ
(2,0) eat	ROOT	(2,0) eat	ROOT
(3,5) the	DET	(3,4) the	DET
(4,5) cheese	MOD	(4,2) cheese	OBJ
(5,2) sandwich	OBJ	(5,2) sandwich	PRED

Accuracy =  $\frac{2}{5}$  = 40%

# Neural parsing

# Neural parsing – simple approach

(Jia and Liang, 2016; Dong and Lapata, 2016)



Some of the very first works doing this.

# Neural parsing – simple approach

- Parsing as translation
  - Linearise grammar from treebank: convert tree to bracketed representation, all in one line
  - Extract sentences from bracketed representation
  - Pair sentences and their linearised trees
  - No need to compute/represent probabilities - **learning**


# Neural parsing – simple approach

- Train sequence to sequence model to *translate* from
  - Sentences to linearised tree; **brackets are tokens!**
  - Use, e.g. LSTM, transformers etc...
  - Attention helps
  - Train with cross-entropy
  - Evaluate as a translation (BLEU) or parsing tasks (parseval)

# Neural parsing – advanced approaches

Table by Chris Manning

- Graph-based methods



Method	UAS	LAS (PTB WSJ SD 3.3)
Chen & Manning 2014	92.0	89.7
Weiss et al. 2015	93.99	92.05
Andor et al. 2016	94.61	92.79
<b>Dozat &amp; Manning 2017</b>	<b>95.74</b>	<b>94.08</b>

- UAS: Unlabeled attachment score
- LAS: Labeled attachment score

Not much more progress from there

# Discussion

- Parsing is an important step for many applications
- Statistical models such as PCFGs allow for resolution of ambiguities
  - Lexicalisation and non-terminal splitting are required to effectively better resolve many ambiguities
- Current statistical/neural parsers are quite accurate
  - ~95% dependency; 97% constituency
  - Human-expert agreement: ~98%