

# Dynamic Programming: Theory and Practice

Alexander Rush and Nathaniel Filardo  
MASC 2013

October 11, 2013

So I've been working with a student.

So I've been working with a student.

And I ask her to build a dependency parser.

So I've been working with a student.

And I ask her to build a dependency parser.

“What should I use to write that?”

## A Dynamic Programming Algorithm

- ▶ Base case definition: for all  $i = 1 \dots n$ , for  $X \in N$

$$\pi[i, i, X] = q(X \rightarrow w_i)$$

(note: define  $q(X \rightarrow w_i) = 0$  if  $X \rightarrow w_i$  is not in the grammar)

- ▶ Recursive definition: for all  $i = 1 \dots n$ ,  $j = (i + 1) \dots n$ ,  $X \in N$ ,

$$\pi(i, j, X) = \max_{\substack{X \rightarrow YZ \in R, \\ s \in \{i \dots (j-1)\}}} (q(X \rightarrow YZ) \times \pi(i, s, Y) \times \pi(s + 1, j, Z))$$

when I am in my advisor's class

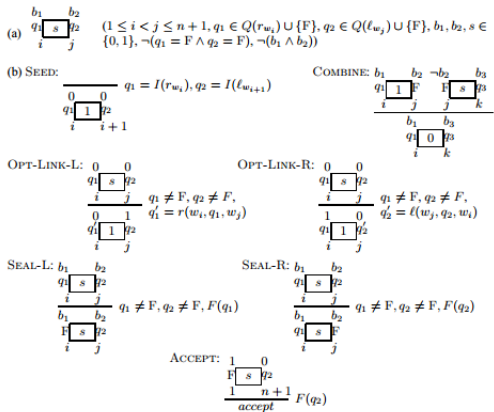


Figure 1.3 Declarative specification of an  $O(n^3)$  algorithm. (a) Form of items in the parse chart. (b) Inference rules. As in Fig. 1.2b,  $F$  is a literal that means ‘an unspecified final state.’

when I read the paper

```

39 // static type for each edge: run time  $O(n^3 + Tn^2)$   $T$  is number of types
40 public Object[][] decodeProjective(DependencyInstance inst,
41     FeatureVector[][] fvs,
42     double[][] probs,
43     FeatureVector[][] nt_fvs,
44     double[][] nt_probs, int K) {
45
46     String[] forms = inst.forms;
47     String[] pos = inst.postags;
48
49     int[][] static_types = null;
50     if(pipe.labeled) {
51         static_types = getTypes(nt_probs, forms.length);
52     }
53
54     KBestParseForest pf = new KBestParseForest(0, forms.length-1, inst, K);
55
56     for(int s = 0; s < forms.length; s++) {
57         pf.add(s, -1, 0, 0.0, new FeatureVector());
58         pf.add(s, -1, 1, 0.0, new FeatureVector());
59     }
60
61     for(int j = 1; j < forms.length; j++) {
62         for(int s = 0; s < forms.length && s+j < forms.length; s++) {
63             int t = s+j;
64
65             FeatureVector prodFV_st = fvs[s][t][0];
66             FeatureVector prodFV_ts = fvs[s][t][1];
67             double prodProb_st = probs[s][t][0];
68             double prodProb_ts = probs[s][t][1];
69
70             int type1 = pipe.labeled ? static_types[s][t] : 0;
71             int type2 = pipe.labeled ? static_types[t][s] : 0;
72
73             FeatureVector nt_fv_s_01 = nt_fvs[s][type1][0][1];
74             FeatureVector nt_fv_s_10 = nt_fvs[s][type2][1][0];
75             FeatureVector nt_fv_t_00 = nt_fvs[t][type1][0][0];
76             FeatureVector nt_fv_t_11 = nt_fvs[t][type2][1][1];
77             double nt_prob_s_01 = nt_probs[s][type1][0][1];
78             double nt_prob_s_10 = nt_probs[s][type2][1][0];
79             double nt_prob_t_00 = nt_probs[t][type1][0][0];
80             double nt_prob_t_11 = nt_probs[t][type2][1][1];
81
82             double prodProb = 0.0;
83
84             for(int r = s; r <= t; r++) {
85
86                 /** first is direction, second is complete*/
87                 /** _s means s is the parent*/
88                 if(r != t) {

```

when I look at the code

## Question

We no longer manually implement

- ▶ Simplex for LPs
- ▶ SVM training for classification
- ▶ FST algorithms
- ▶ etc.



## Question

We no longer manually implement

- ▶ Simplex for LPs
- ▶ SVM training for classification
- ▶ FST algorithms
- ▶ etc.

But I don't have a good answer for dynamic programming.

# Outline

## Four views of Dynamic Programming

1. Imperative k
2. Declarative
3. Graphical
4. Optimization

## Dynamic Programming for NLP

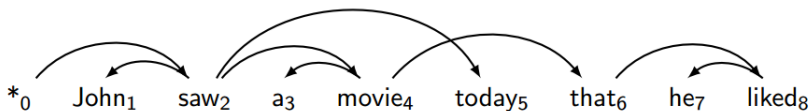
# Decoding for NLP

- ▶  $f$ ; a scoring function.
- ▶  $y$ ; a possible output structure.

**optimization problem:**

$$\max_y f(y)$$

## Example: Dependency Parsing



$$\begin{aligned} f(y) = & \text{score}(\text{head} = *_{0}, \text{mod} = \text{saw}_{2}) + \text{score}(\text{saw}_{2}, \text{John}_{1}) \\ & + \text{score}(\text{saw}_{2}, \text{movie}_{4}) + \text{score}(\text{saw}_{2}, \text{today}_{5}) \\ & + \text{score}(\text{movie}_{4}, \text{a}_{3}) + \dots \end{aligned}$$

# Dynamic Programming for Decoding

“Dynamic programming is a method for solving complex problems by breaking them down into simpler subproblem”

Traditionally defined as an optimization problem defined using the Bellman equations.

However these equations don't give use much insight into what to do in practice.

View 1: Imperative

In my advisor's class, we teach an imperative implementation of dynamic programming using charts.



In my advisor's class, we teach an imperative implementation of dynamic programming using charts.

Some students enjoy this, it correlates well with not having bugs.

Initialization:  $C[s][s][d][c] = 0.0 \quad \forall s, d, c$

for  $k : 1..n$

  for  $s : 1..n$

$t = s + k$

    if  $t > n$  then break

      % First: create incomplete items

$$C[s][t][\leftarrow][0] = \max_{s \leq r < t} (C[s][r][\rightarrow][1] + C[r+1][t][\leftarrow][1] + s(t, s)) \quad (*)$$

$$C[s][t][\rightarrow][0] = \max_{s \leq r < t} (C[s][r][\rightarrow][1] + C[r+1][t][\leftarrow][1] + s(s, t))$$

      % Second: create complete items

$$C[s][t][\leftarrow][1] = \max_{s \leq r < t} (C[s][r][\leftarrow][1] + C[r][t][\leftarrow][0])$$

$$C[s][t][\rightarrow][1] = \max_{s < r \leq t} (C[s][r][\rightarrow][0] + C[r][t][\rightarrow][1])$$

  end for

end for

# Benefits

- ▶ Often fastest method in practice.
- ▶ Doesn't require any additional code or library.
- ▶ Great way to prove how hardcore you are.

## Issue 1

- ▶ Debugging is very subtle.

$$C[s][t][\leftarrow][0] = \max_{s \leq r < t}$$

$$C[s][t][\rightarrow][0] = \max_{s \leq r < t}$$

% Second: create complete

$$C[s][t][\leftarrow][1] = \max_{s \leq r < t}$$

$$C[s][t][\rightarrow][1] = \max_{s < r \leq t}$$



Thanks for your post, I'm also struggling with making all of this work :( . Here is what I have done so far. I have taken the test sentence:

```

STAT5A
mutations
in
the
Src
homology
2
{
SH2
}
and
SH3
domains
did
not
alter
the
BTK
-
mediated
tyrosine
phosphorylation
.

```

Sentence has 23 words so first one is at  $k = 0$  and last one at  $k = 22$  (it being the . character)

And when I run my code I get the following debug output:

```

k = 0 U = 0 V = 0
Calculating  $P_i[0, *, *] * q(0|*, *) * e(\text{STAT5A} \mid I\_GENE)$ 
Taken max probability = 0.011542410287450856
k = 0 U = 0 V = 1
Calculating  $P_i[0, *, *] * q(1|*, *) * e(\text{STAT5A} \mid O)$ 
Taken max probability = 0.004527456817631497
k = 0 U = 1 V = 0
Calculating  $P_i[0, *, *] * q(0|*, *) * e(\text{STAT5A} \mid I\_GENE)$ 
Taken max probability = 0.011542410287450856
k = 0 U = 1 V = 1
Calculating  $P_i[0, *, *] * q(1|*, *) * e(\text{STAT5A} \mid O)$ 
Taken max probability = 0.004527456817631497

```

## Issue 2

- ▶ Optimized to the specific problem domain.

Cannot use MSTParser for POS tagging.

## Issue 2

- ▶ Optimized to the specific problem domain.

Cannot use MSTParser for POS tagging.

Cannot really use MSTParser third-order parsing.

## Issue 3

Cannot easily utilize dynamic programming extensions.

- ▶ Variant algorithms
  - ▶ K-Best
  - ▶ Loss-Augmented inference
  - ▶ Posterior Decoding



## Issue 3

Cannot easily utilize dynamic programming extensions.

- ▶ Variant algorithms
  - ▶ K-Best
  - ▶ Loss-Augmented inference
  - ▶ Posterior Decoding
- ▶ Pruning algorithms
  - ▶ Max-Marginals
  - ▶ Coarse-to-Fine

## Issue 3

Cannot easily utilize dynamic programming extensions.

- ▶ Variant algorithms
  - ▶ K-Best
  - ▶ Loss-Augmented inference
  - ▶ Posterior Decoding
- ▶ Pruning algorithms
  - ▶ Max-Marginals
  - ▶ Coarse-to-Fine
- ▶ Constrained algorithms
  - ▶ Beam Search
  - ▶ Dual Decomposition
  - ▶ Lagrangian Relaxation

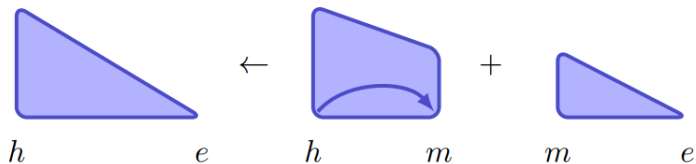
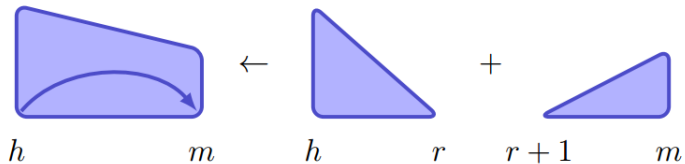
....

## Opinion Slide

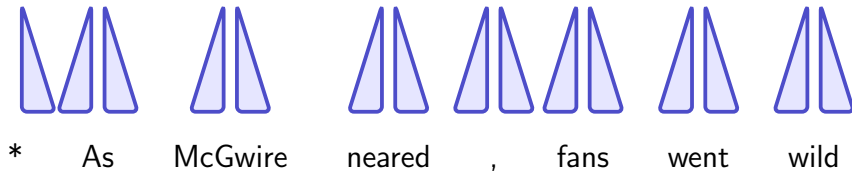
Do not write imperative dynamic programming code.

## View 2: Declarative

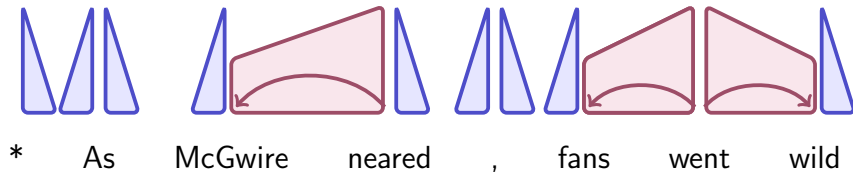
# Dependency Parsing Setup



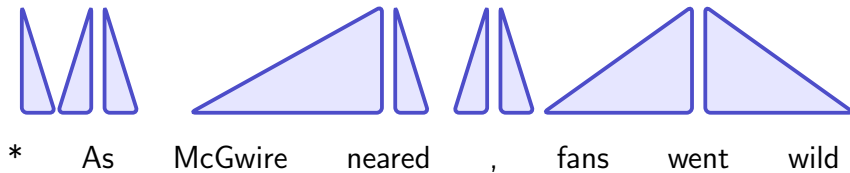
# First-Order Parsing



# First-Order Parsing

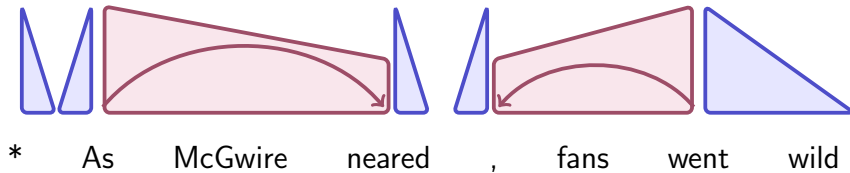


# First-Order Parsing

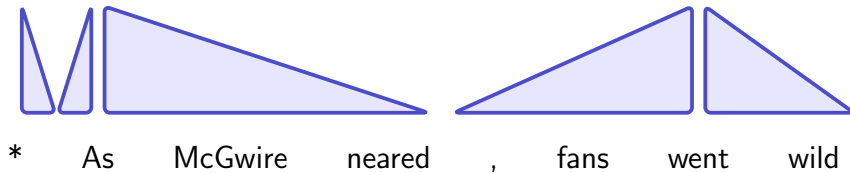




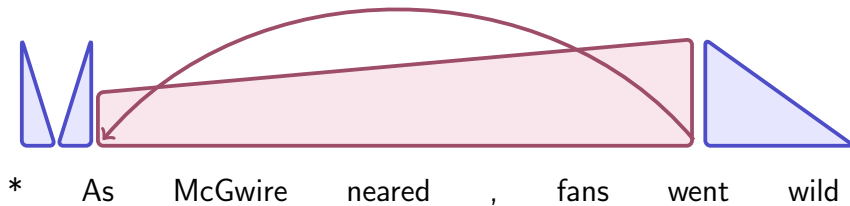
# First-Order Parsing



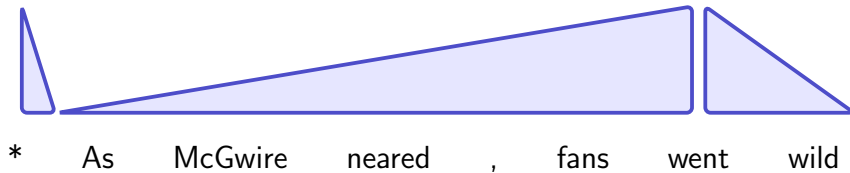
# First-Order Parsing



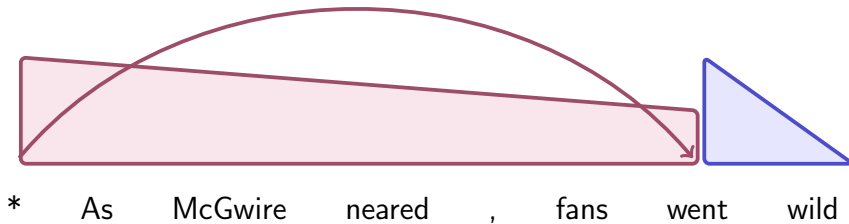
# First-Order Parsing



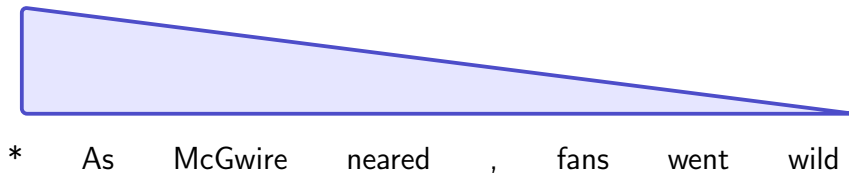
# First-Order Parsing



# First-Order Parsing



## First-Order Parsing



Declarative code is clean, often less buggy.

```
:- item(item, double, 0).
:- item([lhalf,rhalf,end], bool, false).
constit(A,H,H,H) += rewrite(A:D,D) if word(D,H).
left(A/C,I,J,H) += constit(B,I,H2,J) * rewrite(A:D,B:D2,C:D)
if lhalf(C,=J+1,H) & word(D,H) & word(D2,H2).
right(A\C,H,I,J) += constit(B,I,H2,J) * rewrite(A:D,C:D,B:D2)
if rhalf(C,H,=J-1) & word(D,H) & word(D2,H2).
lhalf(A,I,H) |= constit(A,I,H,J)!=0.
rhalf(A,H,J) |= constit(A,I,H,J)!=0.
constit(A,I,H,K) += left(A/C,I,J,H) * constit(C,=J+1,H,K).
constit(A,I,H,K) += constit(C,I,H,=J-1) * right(A\C,H,J,K).
goal += constit(s,1,_,N) if end(N).
```

Declarative code is clean, often less buggy.

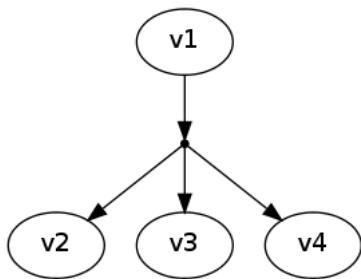
```
:- item(item, double, 0).  
:- item([lhalf,rhalf,end], bool, false).  
constit(A,H,H,H) += rewrite(A:D,D) if word(D,H).  
left(A/C,I,J,H) += constit(B,I,H2,J) * rewrite(A:D,B:D2,C:D)  
if lhalf(C,=J+1,H) & word(D,H) & word(D2,H2).  
right(A\C,H,I,J) += constit(B,I,H2,J) * rewrite(A:D,C:D,B:D2)  
if rhalf(C,H,=J-1) & word(D,H) & word(D2,H2).  
lhalf(A,I,H) |= constit(A,I,H,J)!=0.  
rhalf(A,H,J) |= constit(A,I,H,J)!=0.  
constit(A,I,H,K) += left(A/C,I,J,H) * constit(C,=J+1,H,K).  
constit(A,I,H,K) += constit(C,I,H,=J-1) * right(A\C,H,J,K).  
goal += constit(s,1,_,N) if end(N).
```

I am not smart enough to make this efficient.



## View 3: Hypergraph

# Hyperedge

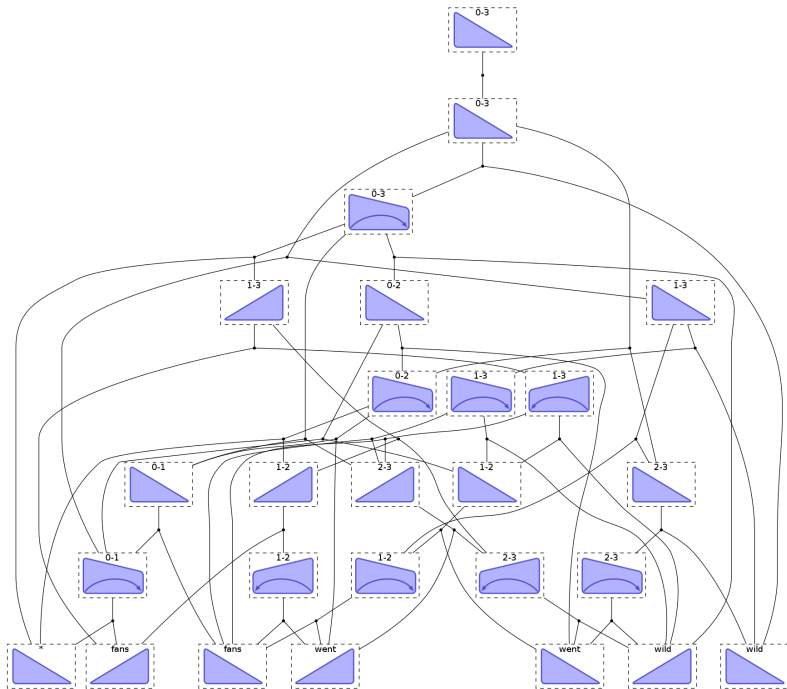


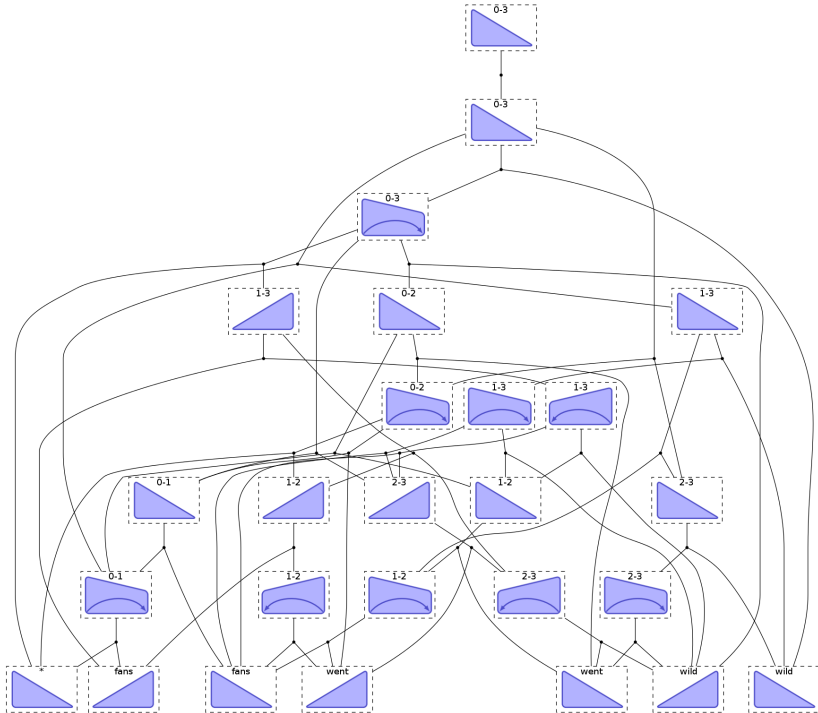
# Hypergraph

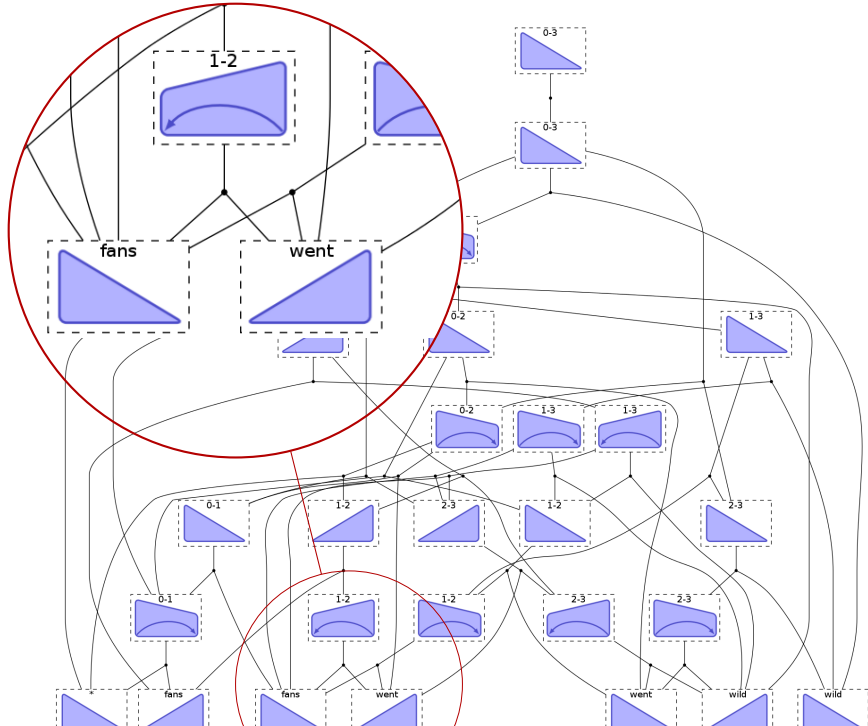
$\mathcal{V}$	set of vertices	1	root node
$\mathcal{E}$	set of hyperedges	$h(e)$	head vertex of edge
		$t(e)$	tail vertex of edge

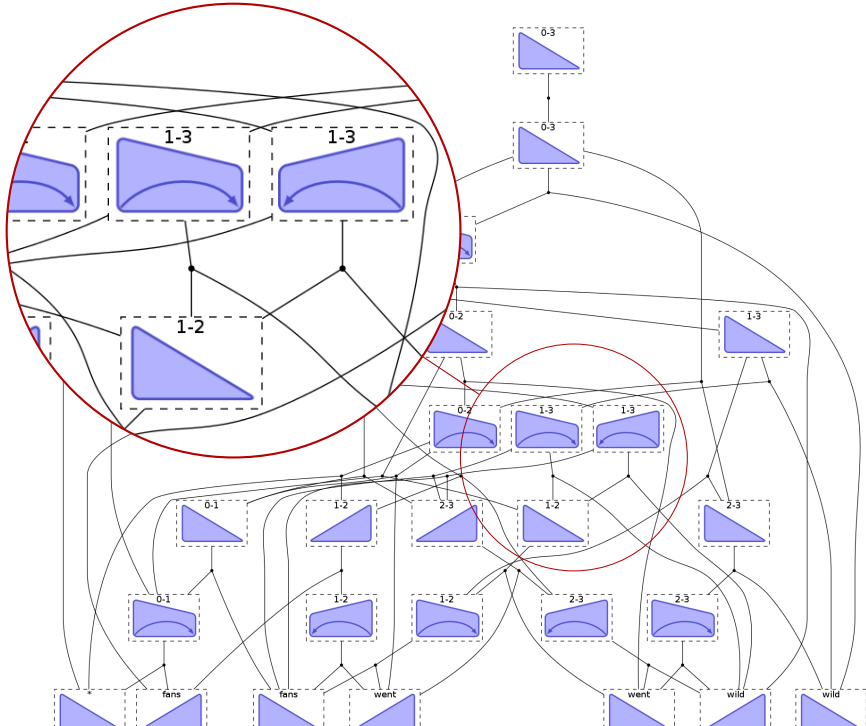
- ▶  $\mathcal{X}$  - set of valid “hyperpaths”  $y$ 
  - ▶  $y \in \{0, 1\}^{|\mathcal{V}|+|\mathcal{E}|}$  has  $y(v) = 1$  iff node  $v$  is in path

Any dynamic program can be represented as a hypergraph.









# Weights

- ▶  $\theta \in R^{|\mathcal{E}|}$  - the weights associated with each hyperedge.

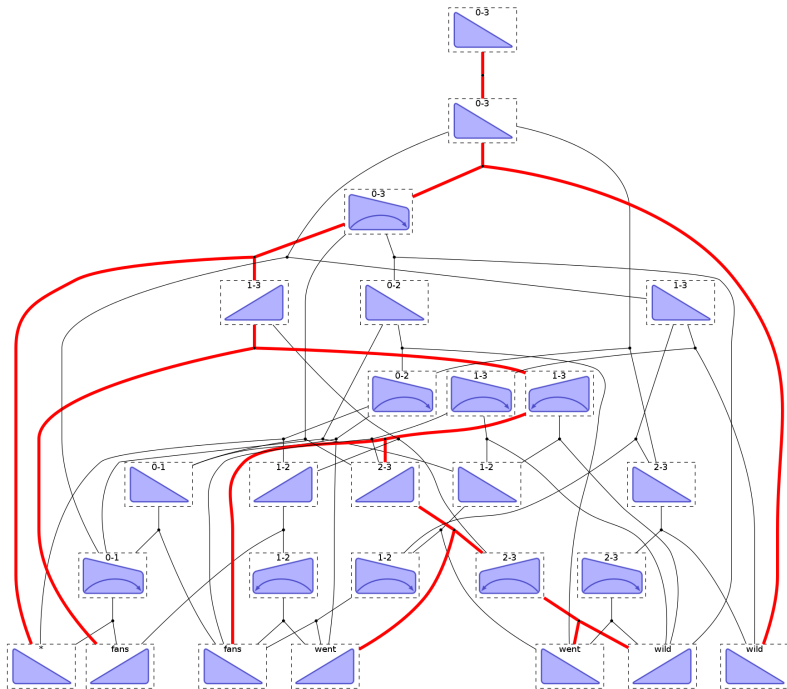
**Hypergraph problem:**

$$\max_{y \in \mathcal{X}} \theta^\top y$$



# Dynamic Programming on Hypergraphs

```
procedure BESTPATH( $\mathcal{V}, \mathcal{E}, \theta$ )  
   $\pi(v) \leftarrow 0$  for all terminals  
  for  $v \in \mathcal{V}$  in bottom-up order do  
     $\pi(v) \leftarrow \max_{e \in \mathcal{E}: h(e)=v} \theta(e) + \pi(t(e))$   
  return  $\pi(1)$ 
```



**procedure** OUTSIDE( $\mathcal{V}, \mathcal{E}, \theta$ )

$\rho[1] \leftarrow 0$

**for**  $e \in \mathcal{E}$  in top-down order **do**

$\langle \langle v_2, \dots, v_k \rangle, v_1 \rangle \leftarrow e$

$t \leftarrow \theta(e) + \sum_{i=2}^k \pi[v_i]$

**for**  $i = 2 \dots k$  **do**

$s \leftarrow \rho[v_1] + t - \pi[v_i]$

**if**  $s > \rho[v_i]$  **then**  $\rho[v_i] \leftarrow s$

**return**  $\rho$

# Benefits

- ▶ Helps debugging, visualizing, pruning, and explaining.
- ▶ Can optimize the hell out of the inner loop code.
- ▶ Useful for other downstream tasks.

# Caveats

- ▶ Need to create edges in memory.
- ▶ Requires construction, often still imperative.

## View 4: Optimization

# Linear Program

- ▶ Any dynamic program can be expressed as a linear program.
- ▶ Can be derived from the Bellman equations.
- ▶ Also directly from the hypergraph representation.

## Linear Constraints Conversion

$$\begin{aligned}\mathcal{X} &= \{y : y(1) = 1, \\ y(v) &= \sum_{e:h(e)=v} y(e) \ \forall v \in \mathcal{V} \text{ except root}, \\ y(v) &= \sum_{e:t(e)=v} y(e) \ \forall v \in \mathcal{V} \text{ except terminals}\}\end{aligned}$$



# Linear Constraints Conversion

$$\mathcal{X} = \{y : y(1) = 1,$$

$$y(v) = \sum_{e:h(e)=v} y(e) \quad \forall v \in \mathcal{V} \text{ except root},$$

$$y(v) = \sum_{e:t(e)=v} y(e) \quad \forall v \in \mathcal{V} \text{ except terminals}\}$$

- $O(|\mathcal{V}|)$  constraints

# Linear Program

$$\max_y \theta^\top y$$

$$y(1) = 1$$

$$y(v) = \sum_{e:h(e)=v} y(e) \quad \forall v \in \mathcal{V} \text{ except root}$$

$$y(v) = \sum_{e:t(e)=v} y(e) \quad \forall v \in \mathcal{V} \text{ except terminals}$$

$$0 \leq y(e) \leq 1 \quad \forall e \in \mathcal{E}$$

$$0 \leq y(v) \leq 1 \quad \forall v \in \mathcal{V}$$

\\* Hypergraph Problem \*\

Minimize

OBJ: 0.865309927772 edge\_0 + 0.805027827013 edge\_1 + 0.719704686  
+ 0.824844977148 edge\_3 + 0.668153201232 edge\_4 + 0.93283382422  
+ 0.551267246091 edge\_6

Subject To

\_C1: node\_13\_tri\_right\_0\_3 = 1

\_C10: - edge\_0 + node\_1\_tri\_right\_1\_1 = 0

\_C11: - edge\_1 + node\_2\_tri\_left\_1\_1 = 0

\_C12: - edge\_2 + node\_3\_tri\_right\_2\_2 = 0

\_C13: - edge\_0 + node\_4\_tri\_left\_2\_2 = 0

\_C14: - edge\_3 + node\_5\_tri\_right\_3\_3 = 0

\_C15: - edge\_2 + node\_6\_tri\_left\_3\_3 = 0

\_C16: - edge\_1 + node\_7\_trap\_left\_1\_2 = 0

\_C17: - edge\_4 + node\_8\_tri\_left\_1\_2 = 0

\_C18: - edge\_3 + node\_9\_trap\_right\_2\_3 = 0

...

# Benefits

- ▶ Can solve with general-purpose LP solver. (Gurobi)
- ▶ There are tricks for outside, K-Best, etc.
- ▶ Real benefit: Adding additional constraints.

# Constrained Paths

**problem:** constrain maximization to valid paths

- ▶  $A \in |b| \times |\mathcal{E}|$ ; a matrix of linear constraints
- ▶  $b \in |b|$ ; a constraint vector

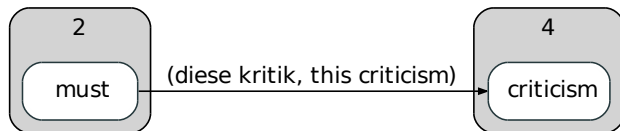
Constrained paths

$$\mathcal{X}' = \{y \in \mathcal{X} : Ay = b\}$$

Constrained best path

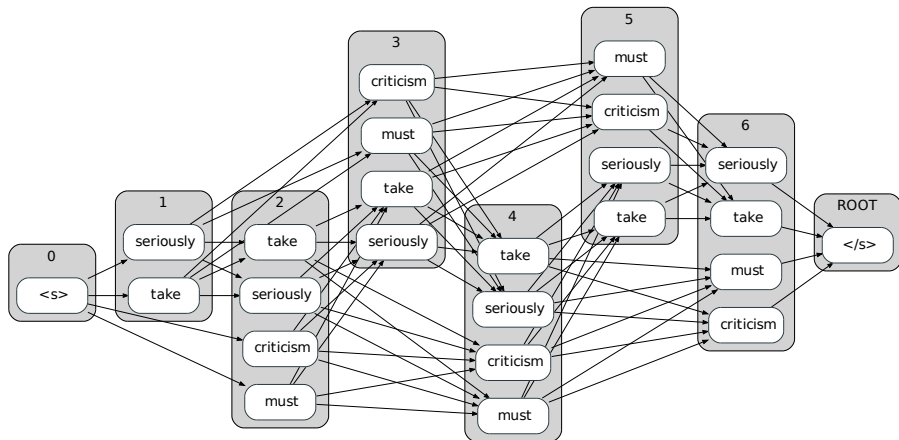
$$\max_{y \in \mathcal{X}'} \theta^\top y$$

## Individual Edge (e)

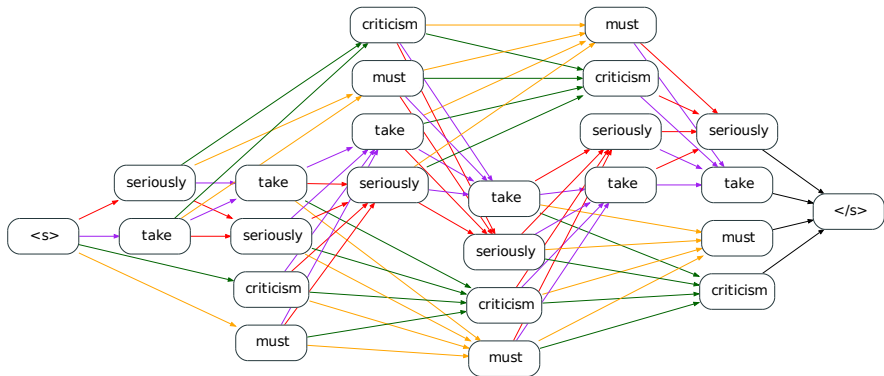


$$\theta(e) = \text{score}(\text{diese kritik}, \text{this criticism}) + \text{score}(\text{must}, \text{this}) + \text{score}(\text{this}, \text{criticism})$$

# Constraints



## Constraints







# Pitch: PyDecode

A pragmatic library for building dynamic programming.

- ▶ Easy interface in Python
- ▶ Hypergraph code in C++
  - ▶ Inside/Outside Viterbi
  - ▶ Pruning with max-marginals
  - ▶ Constraint specification
  - ▶ Dual Decomposition/Lagrangian Relaxation
  - ▶ Export to LP
- ▶ Visualization tools supporting IPython Notebook