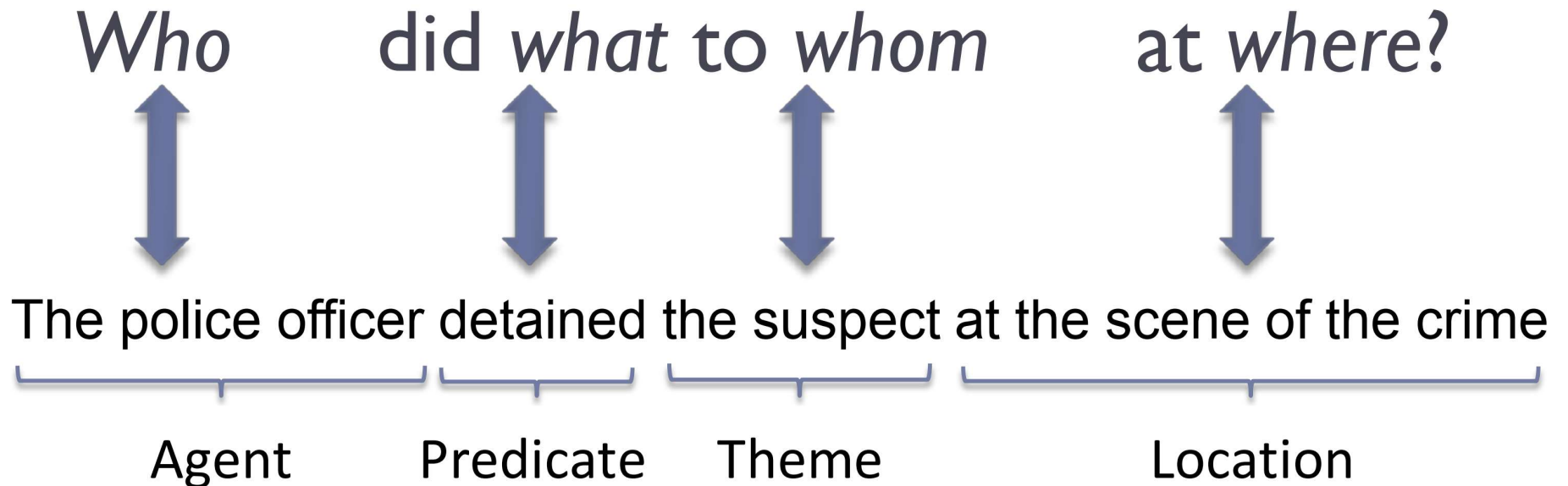


# Semantic Role Labeling

Chong Xiang and Zhongqiao Gao

02/20/2020

# Semantic Role Labeling



B-arg0; O; B-arg1; I-arg0; B-arg2; I-arg2; TMP
B_arg0, B_v, B_arg2, B_arg2, B_arg1, B_arg1, B_argm_temp
B O B I O O O
y1 = B, y2 = O, y3 = B, y4 = I, y5 = B, y6 = I, y7 = O
B_ARG1 O B_ARG2 I_ARG2 O B_ARG3 O
Alice, gave, Bob's, mom, a, book, yesterday
B_ARG0, B_V, B_ARG1, I_ARG1, B_ARG2, I_ARG2, B_ARG3
y1=B-ARG0, y2=B-v, y3=B-ARG1, y4=I-ARG1, y5=B-ARG2, y6=I-ARG2, y7=O
B O B I O B B
B, O, B, I, O, B, B
I'm not <i>*quite*</i> sure, but:
y1=B-ARG0
y2=I-ARG0
y3=B-ARG1
y4=I-ARG1
y5=O
y6=B-ARG2
y7=O
B-ARG1, V, O, I-ARG2, O, I-ARG3, O; I am very confused by the BIO sequence
B_{arg0} B_v B_{arg1} I_{arg1} B_{arg2} I_{arg2} O
lr, O,lr, lr, O,O,O
y1 = B, y2 = O, y3 = B, y4 = I, y5 = B, y6 = I, y7 = O
B(arg0), B(v), B(arg1), I(arg1), B(arg2), I(arg2), O
B, O, B, I, B, I, B
y1=Alice, y2=Bob's mom, y3=a book, modifier=yesterday

# The Proposition Bank (PropBank)

## Give

	ArgM-TMP: when?
	ArgM-LOC: where?
Arg0: giver	ArgM-DIR: where to/from?
Arg1: thing given	ArgM-MNR: how?
Arg2: entity given to	ArgM-PRP/CAU: why?

Alice gave Bob's mom a book yesterday



# The Proposition Bank (PropBank)

## Give

ArgM-TMP: when?

ArgM-LOC: where?

Arg0: giver

ArgM-DIR: where to/from?

Arg1: thing given

ArgM-MNR: how?

Arg2: entity given to

ArgM-PRP/CAU: why?

[Arg0: Alice] gave [Arg2: Bob's mom] [Arg1: a book] [ArgM-TMP: yesterday]

# BIO (Beginning-Inside-Outside) tagging

Alex is going to Los Angeles

**Alex:** B-PER

**is:** O

**going:** O

**to:** O

**Los:** B-LOC

**Angeles:** I-LOC

[Arg0: Alice] gave [Arg2: Bob's mom] [Arg1: a book]  
[ArgM-TMP: yesterday]

**Alice:** B-Arg0

**gave:** B-v

**Bob':** B-Arg2

**mom:** I-Arg2

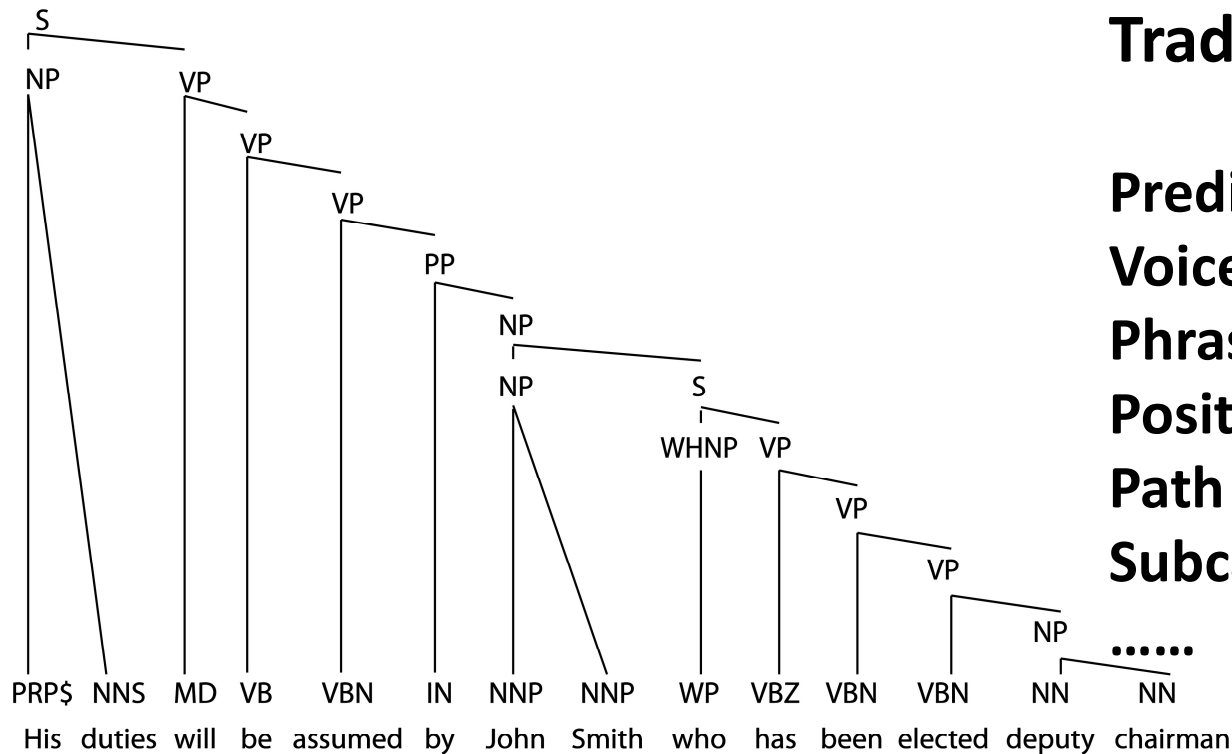
**a:** B-Arg1

**book:** I-Arg1

**Yesterday:** B-ArgM-TMP

# CoNLL-2005

WORDS----->	NE--->	POS	PARTIAL_SYNT	FULL_SYNT----->	VS	TARGETS	PROPS----->		
The	*	DT	(NP*	(S*	(S(NP*	-	-	(A0*	(A0*
\$	*	\$	*	*	(ADJP(QP*	-	-	*	*
1.4	*	CD	*	*	*	-	-	*	*
billion	*	CD	*	*	*)	-	-	*	*
robot	*	NN	*	*	*	-	-	*	*
spacecraft	*	NN	*)	*	*)	-	-	*)	*)
faces	*	VBZ	(VP*	*	(VP*	01	face	(V*	*
a	*	DT	(NP*	*	(NP*	-	-	(A1*	*
six-year	*	JJ	*	*	*	-	-	*	*
journey	*	NN	*)	*	*	-	-	*	*
to	*	TO	(VP*	(S*	(S(VP*	-	-	*	*
explore	*	VB	*)	*	(VP*	01	explore	*	(V*
Jupiter	(ORG*)	NNP	(NP*)	*	(NP(NP*)	-	-	*	(A1*
and	*	CC	*	*	*	-	-	*	*
its	*	PRP\$	(NP*	*	(NP*	-	-	*	*
16	*	CD	*	*	*	-	-	*	*
known	*	JJ	*	*	*	-	-	*	*
moons	*	NNS	*)	*)	*)	-	-	*)	*)
.	*	.	*	*)	*)	-	-	*	*



Clauses:

---

Chunks:

---

NP VP PP NP NP VP NP

Predicate-Argument Structure:

*assume*

---

A1 AM-MOD V A0

---

A1 R-A1 V A2

*elect*

## Traditional features:

**Predicate and POS tag of predicate**

**Voice**

**Phrase type**

**Position**

**Path** S↑NP↑PP↑VP↓VBN

**Subcategorization**

Figure copied from *The Importance of Syntactic Parsing and Inference in Semantic Role Labeling*

- **Feature based Semantic Role Labeling**

Syntax seems to be a prerequisite for SRL

- **Deep Semantic Role Labeling: What Works and What's Next ACL'2017**

End-to-end model without syntactic input

- **Linguistically-Informed Self-Attention for Semantic Role Labeling ACL'2018**

Explicitly model the syntactic information in neural network

- **Feature based Semantic Role Labeling**

Syntax seems to be a prerequisite for SRL

- **Deep Semantic Role Labeling: What Works and What's Next ACL'2017**

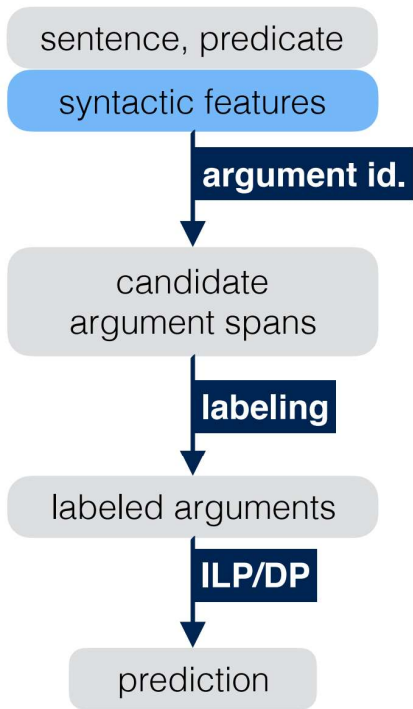
End-to-end model without syntactic input

- **Linguistically-Informed Self-Attention for Semantic Role Labeling ACL'2018**

Explicitly model the syntactic information in neural network

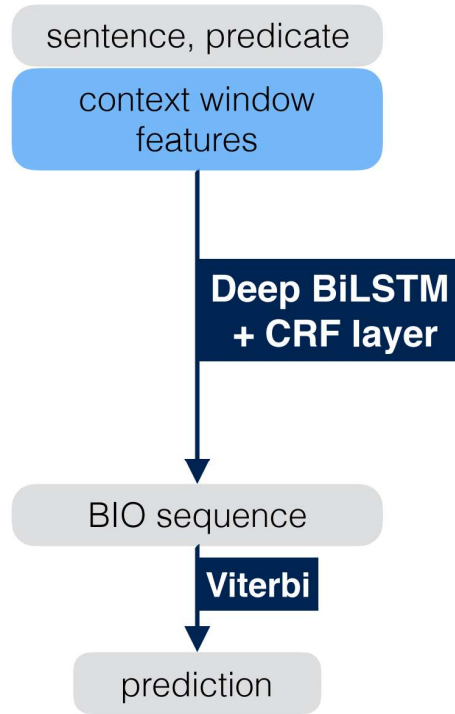
# SRL Systems

## Pipeline Systems



Punyakank et al., 2008  
 Täckström et al., 2015  
 FitzGerald et al., 2015

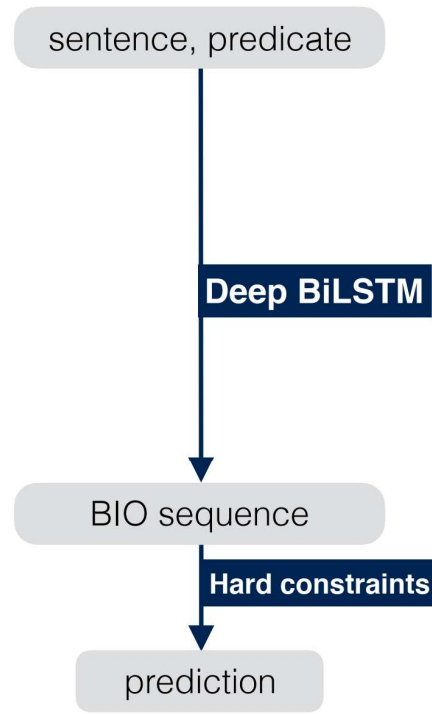
## End-to-end Systems



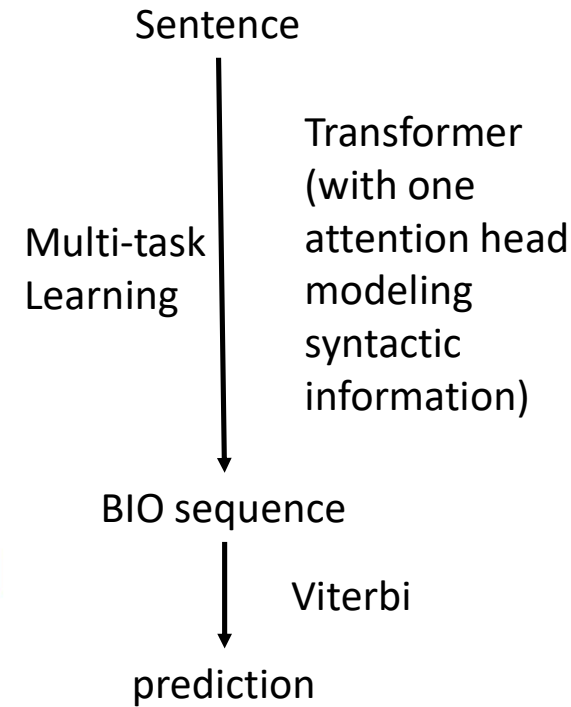
Collobert et al., 2011  
 Zhou and Xu, 2015  
 Wang et. al, 2015

## He et al, 2017

### \*This work



## Strubell et al, 2018



- **Feature based Semantic Role Labeling**

Syntax seems to be a prerequisite for SRL

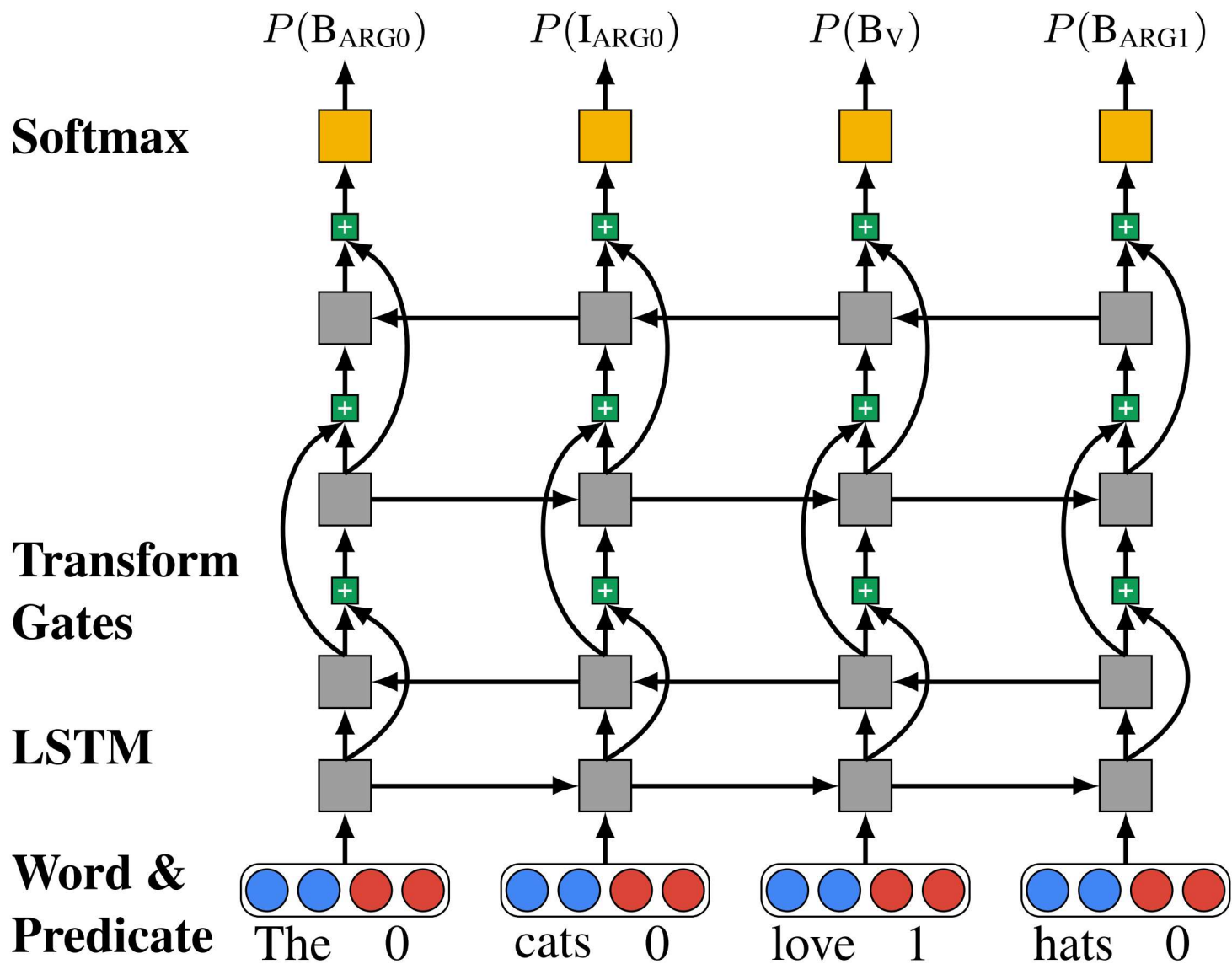
- **Deep Semantic Role Labeling: What Works and What's Next ACL'2017**

End-to-end model without syntactic input

- **Linguistically-Informed Self-Attention for Semantic Role Labeling ACL'2018**

Explicitly model the syntactic information in neural network





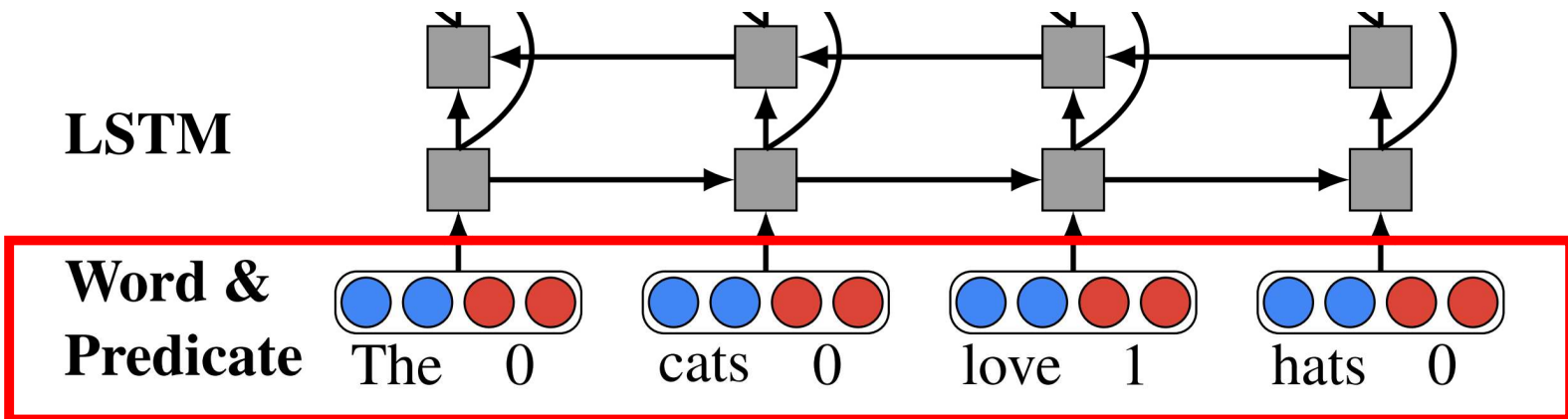
Basic idea: BiLSTM

Design detail?  
Training Technique?

# Input: Glove word embedding and binary predicate mask

A simplification of ACL'2015 (word, predicate, predicate context, and region mark)

$$\mathbf{x}_{l,t} = \begin{cases} [\mathbf{W}_{\text{emb}}(w_t), \mathbf{W}_{\text{mask}}(t = v)] & l = 1 \\ \mathbf{h}_{l-1,t} & l > 1 \end{cases}$$



## Basic structure: BiLSTM

$$i_{l,t} = \sigma(\mathbf{W}_i^l[\mathbf{h}_{l,t+\delta_l}, \mathbf{x}_{l,t}] + \mathbf{b}_i^l)$$

$$o_{l,t} = \sigma(\mathbf{W}_o^l[\mathbf{h}_{l,t+\delta_l}, \mathbf{x}_{l,t}] + \mathbf{b}_o^l)$$

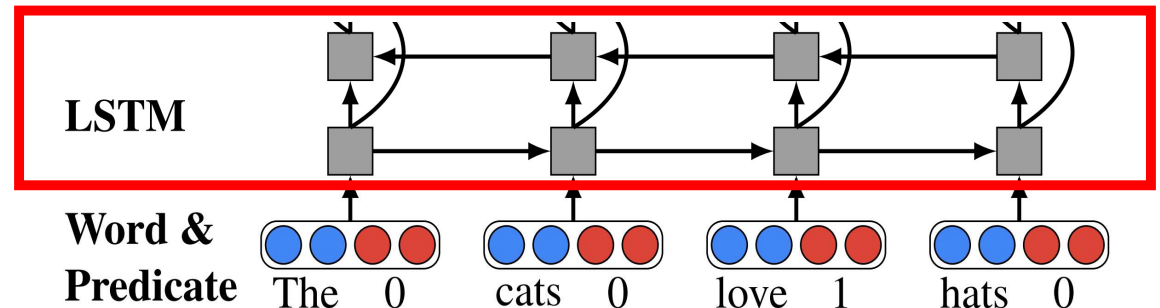
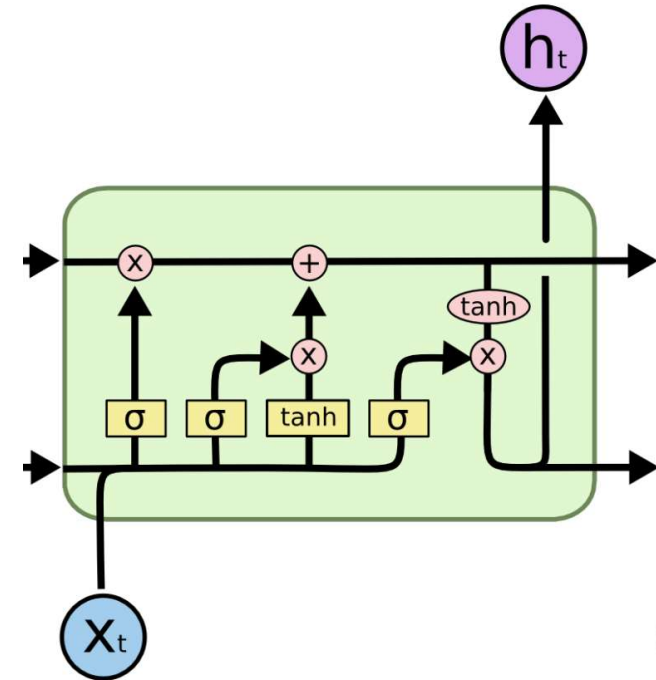
$$f_{l,t} = \sigma(\mathbf{W}_f^l[\mathbf{h}_{l,t+\delta_l}, \mathbf{x}_{l,t}] + \mathbf{b}_f^l + 1)$$

$$\tilde{c}_{l,t} = \tanh(\mathbf{W}_c^l[\mathbf{h}_{l,t+\delta_l}, \mathbf{x}_{l,t}] + \mathbf{b}_c^l)$$

$$\mathbf{c}_{l,t} = i_{l,t} \circ \tilde{c}_{l,t} + f_{l,t} \circ \mathbf{c}_{t+\delta_l}$$

$$\mathbf{h}_{l,t} = o_{l,t} \circ \tanh(\mathbf{c}_{l,t})$$

$$\delta_l = \begin{cases} 1 & \text{if } l \text{ is even} \\ -1 & \text{otherwise} \end{cases}$$



# Network Training: Highway Connections and Recurrent Dropout

**Highway Connection: somewhat like residual network**

$$\mathbf{r}_{l,t} = \sigma(\mathbf{W}_r^l[\mathbf{h}_{l,t-1}, \mathbf{x}_t] + \mathbf{b}_r^l)$$

$$\mathbf{h}'_{l,t} = \mathbf{o}_{l,t} \circ \tanh(\mathbf{c}_{l,t})$$

$$\mathbf{h}_{l,t} = \mathbf{r}_{l,t} \circ \mathbf{h}'_{l,t} + (1 - \mathbf{r}_{l,t}) \circ \mathbf{W}_h^l \mathbf{x}_{l,t}$$

← Transform gate

← LSTM output without gate

← Gated output

**Recurrent Dropout**

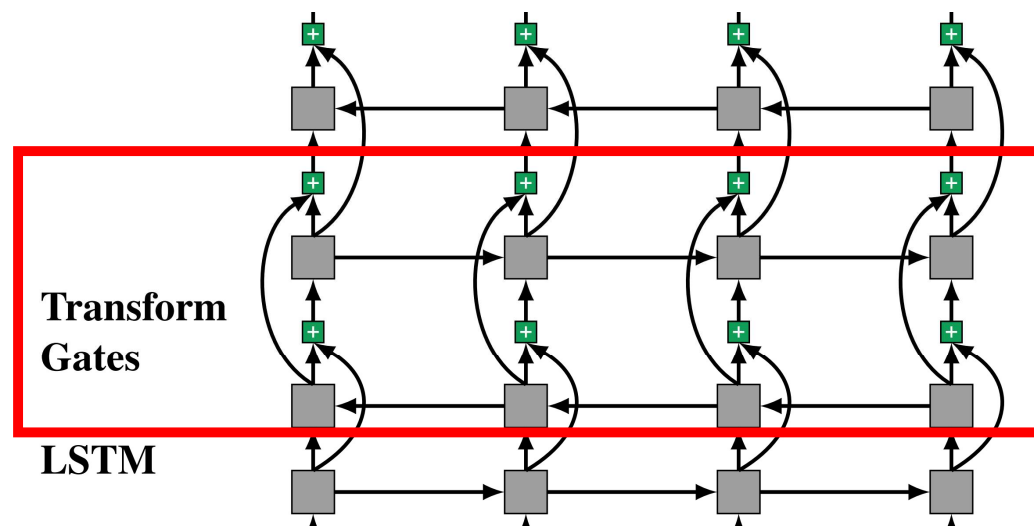
↙ Gated output

$$\tilde{\mathbf{h}}_{l,t} = \mathbf{r}_{l,t} \circ \mathbf{h}'_{l,t} + (1 - \mathbf{r}_{l,t}) \circ \mathbf{W}_h^l \mathbf{x}_{l,t}$$

$$\mathbf{h}_{l,t} = \mathbf{z}_l \circ \tilde{\mathbf{h}}_{l,t}$$

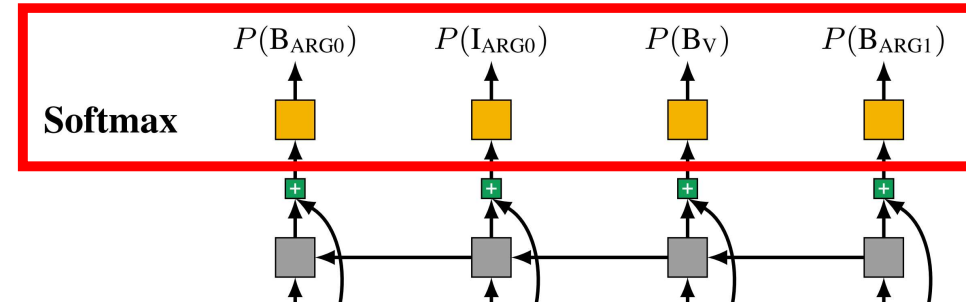
↙ Masked/dropout output

**Orthonormal initialization**



# Output: Constrained A\* decoding

Another difference with ACL'2015 (CRF)



1. Model the dependencies between the output tags
2. Add constraints to reject invalid output
3. Use A\* searching algorithm to find the “optimal” tag sequence

$$p(y_t | \mathbf{x}) \propto \exp(\mathbf{W}_{\text{tag}}^y \mathbf{h}_{L,t} + \mathbf{b}_{\text{tag}})$$

← Softmax output

$$f(\mathbf{w}, y_{1:t}) = \sum_{i=1}^t \log p(y_i | \mathbf{w}) - \sum_{c \in \mathcal{C}} c(\mathbf{w}, y_{1:i})$$

← Confidence value with penalization

$$g(\mathbf{w}, y_{1:t}) = \sum_{i=t+1}^n \max_{y_i \in \mathcal{T}} \log p(y_i | \mathbf{w})$$

← A\* heuristic

# Constraint example

## BIO Constraints

Reject invalid BIO transitions, such as  $B_{\text{ARG0}}$  followed by  $I_{\text{ARG1}}$

## SRL Constraints

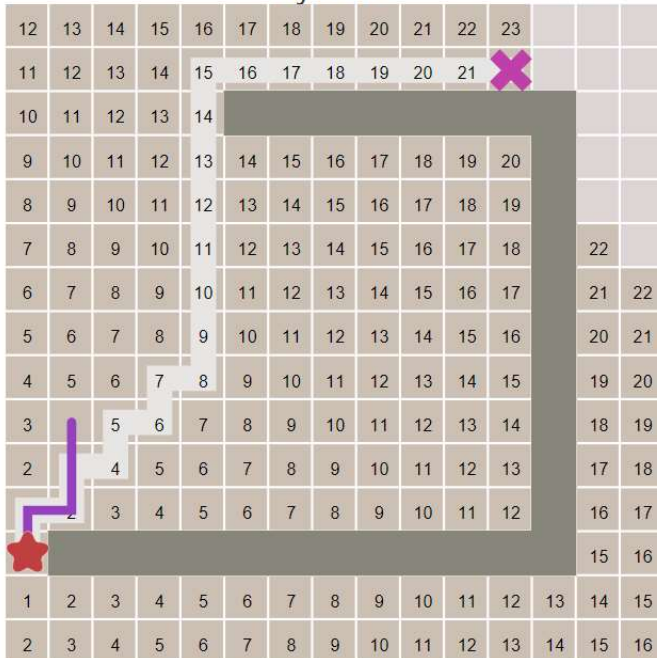
- **Unique core roles (U):** Each core role (ARG0-ARG5, ARG<sub>A</sub>) should appear at most once for each predicate.
- **Continuation roles (C):** A continuation role C-X can exist only when its base role X is realized before it.
- **Reference roles (R):** A reference role R-X can exist only when its base role X is realized (not necessarily before R-X).

## Syntactic Constraints

We can enforce consistency with a given parse tree by rejecting or penalizing arguments that are not constituents.

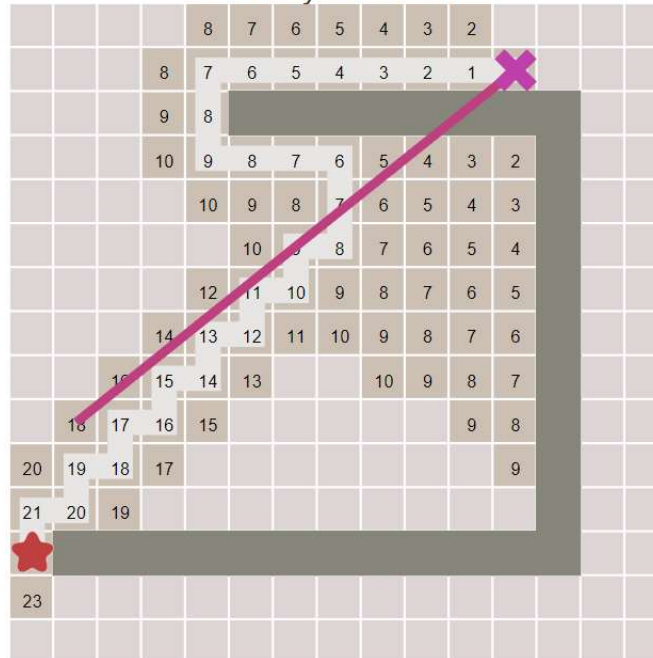
# Intuition on A\* searching

Dijkstra's



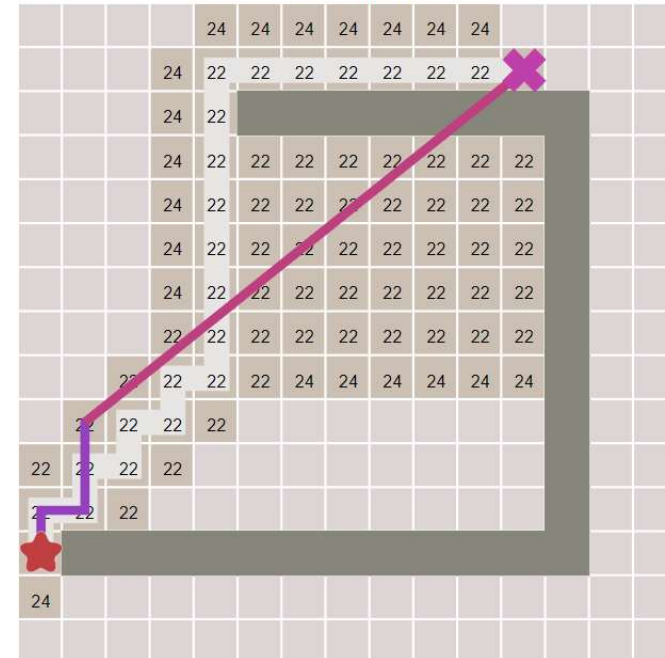
Record the distance from the starting point

Greedy Best-First



Record the distance away from the destination

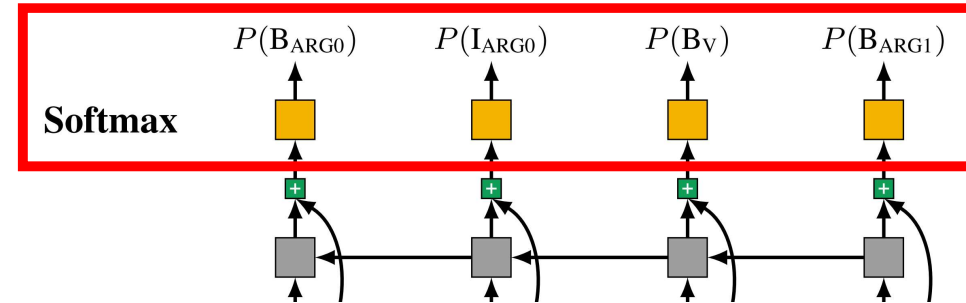
A\* Search



Consider both distances

# Output: Constrained A\* decoding

Another difference with ACL'2015 (CRF)



1. Not used during the training
2. A\* searching algorithm reduced to Viterbi and able find the optimal solution?

$$p(y_t | \mathbf{x}) \propto \exp(\mathbf{W}_{\text{tag}}^y \mathbf{h}_{L,t} + \mathbf{b}_{\text{tag}})$$

← Softmax output

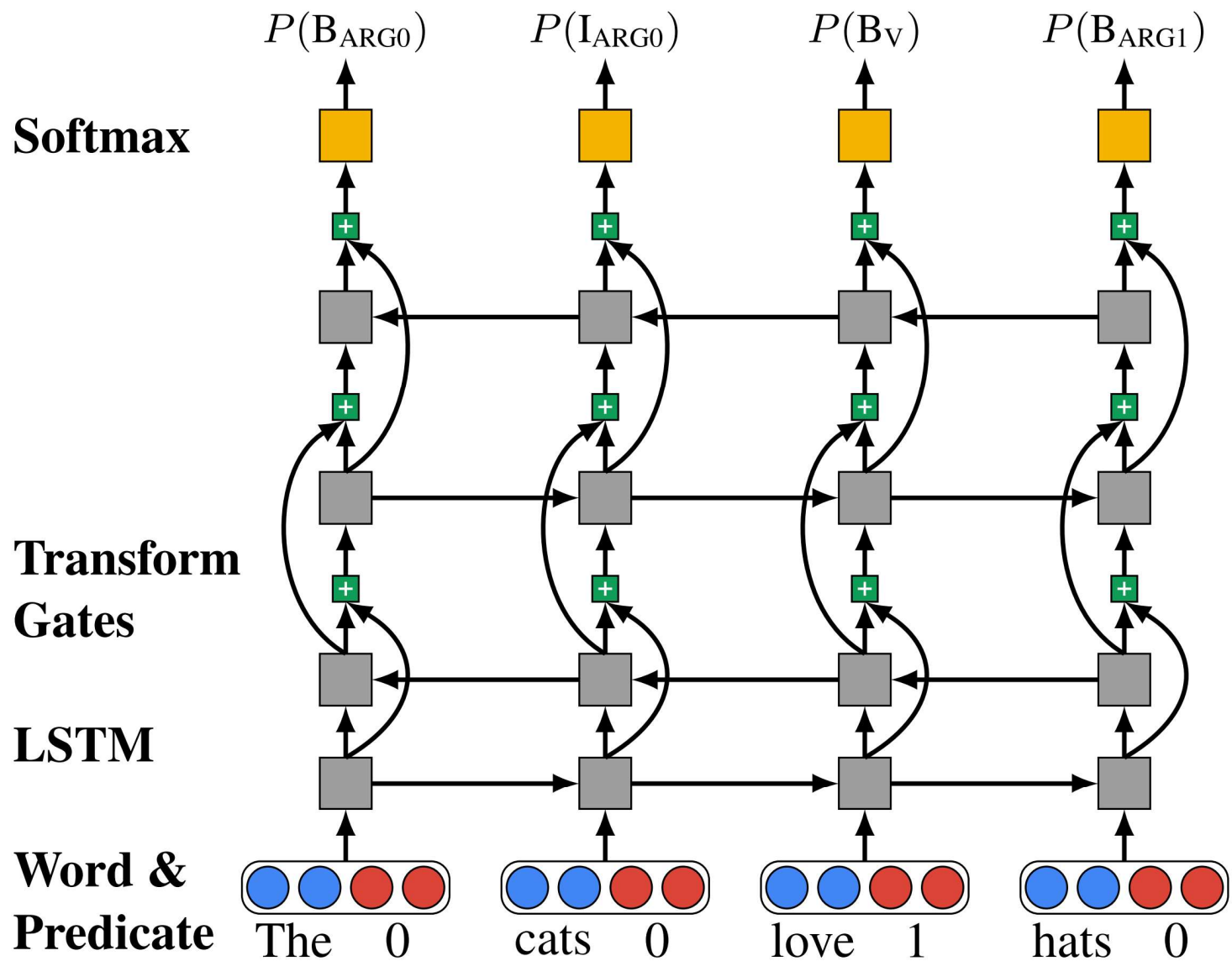
$$f(\mathbf{w}, y_{1:t}) = \sum_{i=1}^t \log p(y_i | \mathbf{w}) - \sum_{c \in \mathcal{C}} c(\mathbf{w}, y_{1:i})$$

← Confidence value with penalization

$$g(\mathbf{w}, y_{1:t}) = \sum_{i=t+1}^n \max_{y_i \in \mathcal{T}} \log p(y_i | \mathbf{w})$$

← A\* heuristic





**CoNLL-2005  
(PropBank)**

**CoNLL-2012  
(OntoNotes)**

Size

40k sentences

140k sentences

Domains

- WSJ / newswire
- Brown (test-only)

- telephone conversations
- newswire
- newsgroups
- broadcast news
- broadcast conversation
- weblogs

Annotated  
predicates

Verbs

Added some nominal  
predicates

“George III is the king of England”

# CoNLL-2005

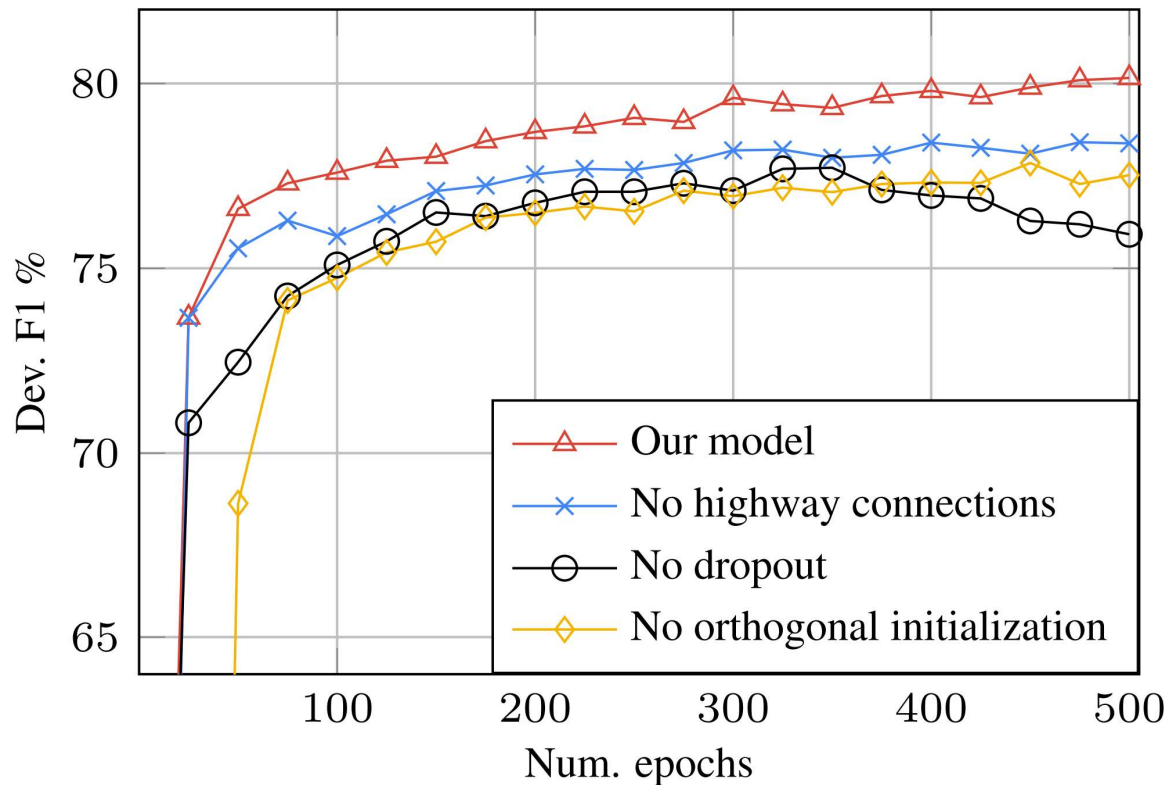
WORDS----->	NE---->	POS	PARTIAL_SYNT	FULL_SYNT----->	VS	TARGETS	PROPS----->
The	*	DT	(NP* (S*	(S(NP*	-	-	(A0* (A0*
\$	*	\$	* *	(ADJP(QP*	-	-	* *
1.4	*	CD	* *	*	-	-	* *
billion	*	CD	* *	*)	-	-	* *
robot	*	NN	* *	*	-	-	* *
spacecraft	*	NN	*) *	*)	-	-	*) *)
faces	*	VBZ	(VP*) *	(VP*	01	face	(V*) *
a	*	DT	(NP*	(NP*	-	-	(A1* *
six-year	*	JJ	* *	*	-	-	* *
journey	*	NN	*) *	*	-	-	* *
to	*	TO	(VP* (S*	(S(VP*	-	-	* *
explore	*	VB	*) *	(VP*	01	explore	* (V*)
Jupiter	(ORG*)	NNP	(NP*) *	(NP(NP*)	-	-	* (A1*
and	*	CC	* *	*	-	-	* *
its	*	PRP\$	(NP*	(NP*	-	-	* *
16	*	CD	* *	*	-	-	* *
known	*	JJ	* *	*	-	-	* *
moons	*	NNS	*) *)	*)*)*)*)*)	-	-	*) *)
.	*	.	* *)	*)	-	-	* *

# Model Performance

Method	Development				WSJ Test				Brown Test				Combined
	P	R	F1	Comp.	P	R	F1	Comp.	P	R	F1	Comp.	F1
Ours (PoE)	<b>83.1</b>	<b>82.4</b>	<b>82.7</b>	<b>64.1</b>	<b>85.0</b>	<b>84.3</b>	<b>84.6</b>	<b>66.5</b>	<b>74.9</b>	<b>72.4</b>	<b>73.6</b>	<b>46.5</b>	<b>83.2</b>
Ours	81.6	81.6	81.6	62.3	83.1	83.0	83.1	64.3	72.9	71.4	72.1	44.8	81.6
Zhou	79.7	79.4	79.6	-	82.9	82.8	82.8	-	70.7	68.2	69.4	-	81.1
FitzGerald (Struct.,PoE)	81.2	76.7	78.9	55.1	82.5	78.2	80.3	57.3	74.5	70.0	72.2	41.3	-
Täckström (Struct.)	81.2	76.2	78.6	54.4	82.3	77.6	79.9	56.0	74.3	68.6	71.3	39.8	-
Toutanova (Ensemble)	-	-	78.6	58.7	81.9	78.8	80.3	60.1	-	-	68.8	40.8	-
Punyakankok (Ensemble)	80.1	74.8	77.4	50.7	82.3	76.8	79.4	53.8	73.4	62.9	67.8	32.3	77.9

Table 1: Experimental results on CoNLL 2005, in terms of precision (P), recall (R), F1 and percentage of completely correct predicates (Comp.). We report results of our best single and ensemble (PoE) model. The comparison models are Zhou and Xu (2015), FitzGerald et al. (2015), Täckström et al. (2015), Toutanova et al. (2008) and Punyakankok et al. (2008).

# Contributions of Three Training Techniques



**1. Without any of the three, seems unable to beat state-of-the-art**

**2. Orthogonal initialization is very important at the early stage of training**

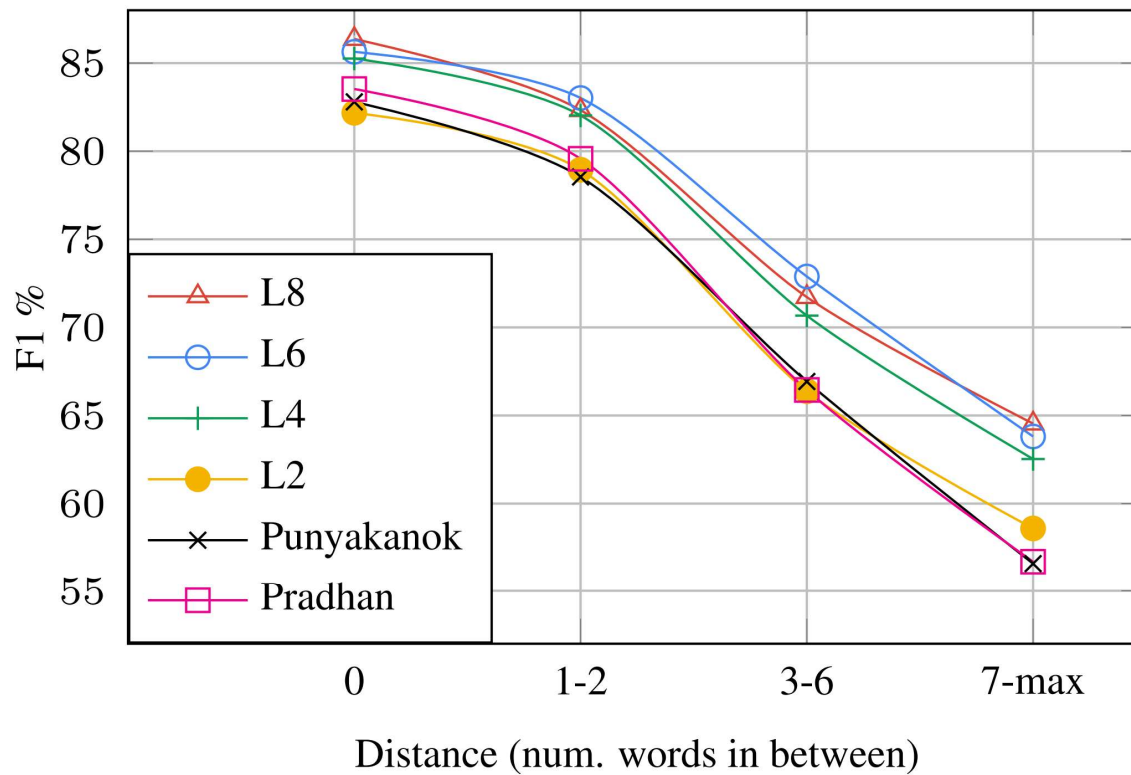
# End-to-End SRL

Train a separate predicate detection model

Dataset	Predicate Detection			End-to-end SRL (Single)			End-to-end SRL (PoE)			$\Delta$ F1
	P	R	F1	P	R	F1	P	R	F1	
CoNLL 2005 Dev.	97.4	97.4	97.4	80.3	80.4	80.3	81.8	81.2	81.5	-1.2
WSJ Test	94.5	98.5	96.4	80.2	82.3	81.2	82.0	83.4	82.7	-1.9
Brown Test	89.3	95.7	92.4	67.6	69.6	68.5	69.7	70.5	70.1	-3.5
CoNLL 2012 Dev.	88.7	90.6	89.7	74.9	76.2	75.5	76.5	77.8	77.2	-6.2
CoNLL 2012 Test	93.7	87.9	90.7	78.6	75.1	76.8	80.2	76.6	78.4	-5.0

Table 3: Predicate detection performance and end-to-end SRL results using predicted predicates.  $\Delta$  F1 shows the absolute performance drop compared to our best ensemble model with gold predicates.

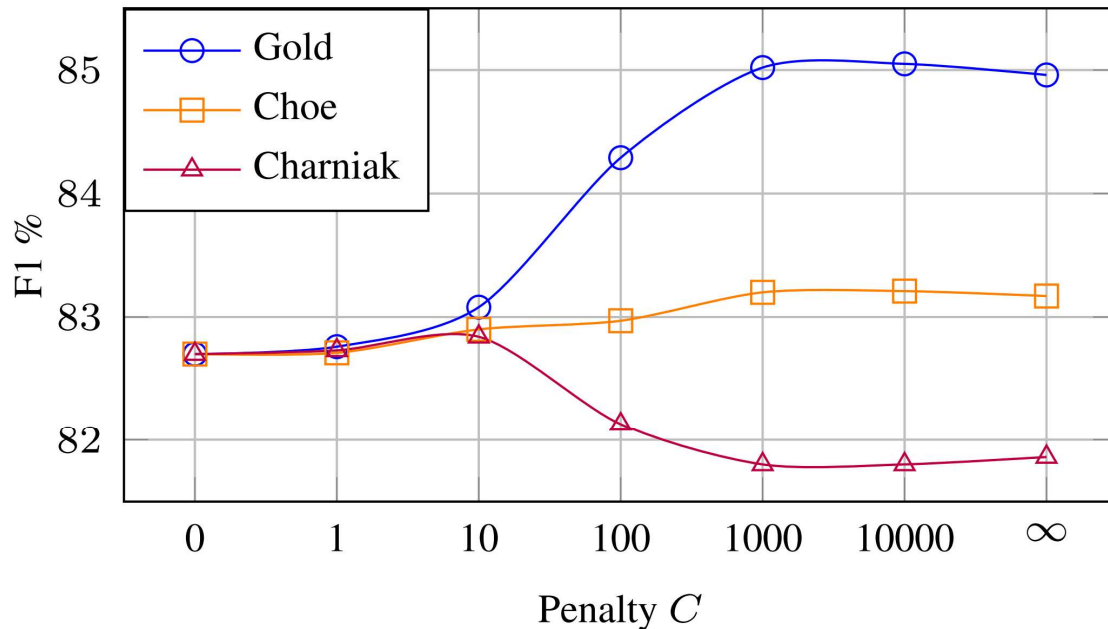
# Long-range Dependency



**Performance deteriorates as distance increases.**

# Adding Syntactic Constraints

$$f(\mathbf{w}, y_{1:t}) = \sum_{i=1}^t \log p(y_i | \mathbf{w}) - \sum_{c \in \mathcal{C}} c(\mathbf{w}, y_{1:i})$$



1. **Syntax helps!**

2. **We need a better automatic parser (syntax modeling), though**

**Gold:** Penn Treebank constituents.

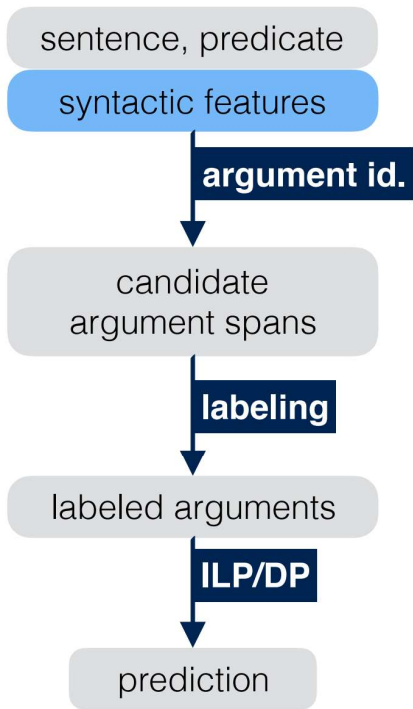
**Choe:** Parsing as language modeling, Choe and Charniak, 2016 (SOTA)

**Charniak:** A maximum-entropy-inspired parser, Charniak, 2000



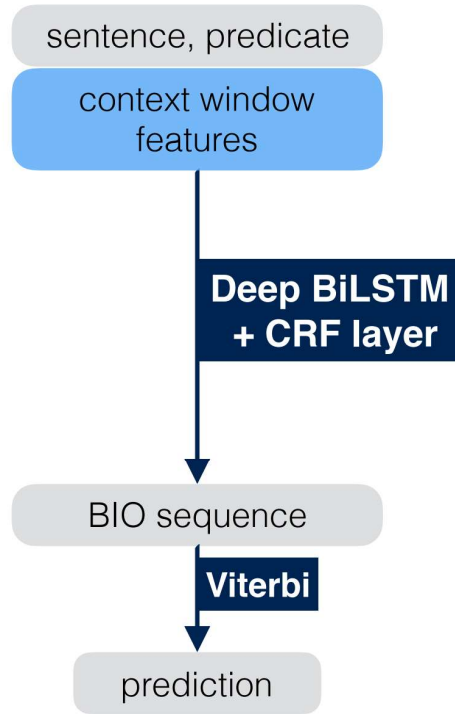
# SRL Systems

## Pipeline Systems



Punyakank et al., 2008  
 Täckström et al., 2015  
 FitzGerald et al., 2015

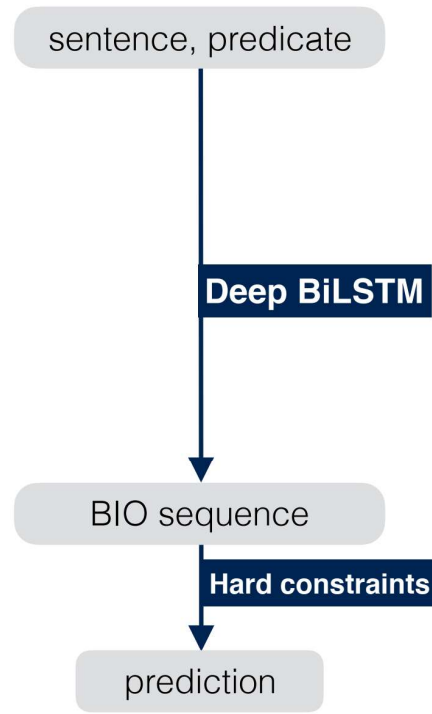
## End-to-end Systems



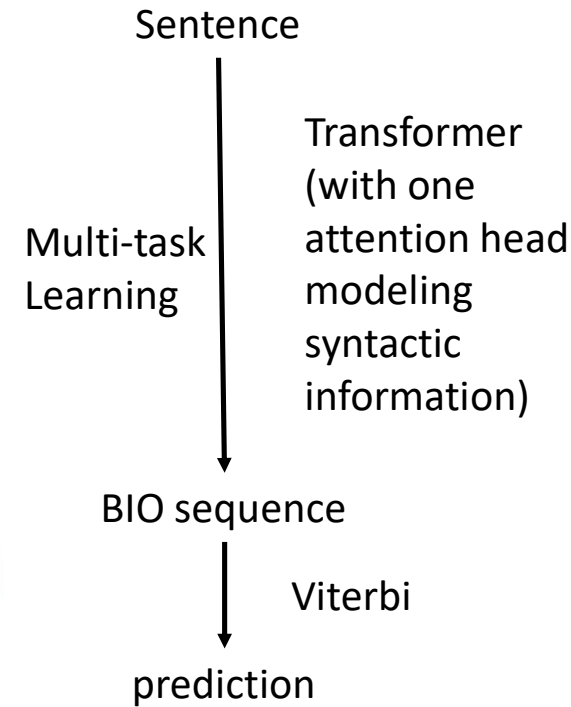
Collobert et al., 2011  
 Zhou and Xu, 2015  
 Wang et. al, 2015

## He et al, 2017

### \*This work



## Strubell et al, 2018



- **Feature based Semantic Role Labeling**

Syntax plays an important role

- **Deep Semantic Role Labeling: What Works and What's Next ACL'2017**

End-to-end model without syntactic input

- **Linguistically-Informed Self-Attention for Semantic Role Labeling ACL'2018**

Explicitly model the syntactic information in neural network

# LISA

## Linguistically-Informed Self-Attention for Semantic Role Labeling



Emma  
Strubell<sup>1</sup>



Patrick  
Verga<sup>1</sup>



Daniel  
Andor<sup>2</sup>



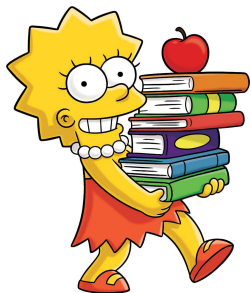
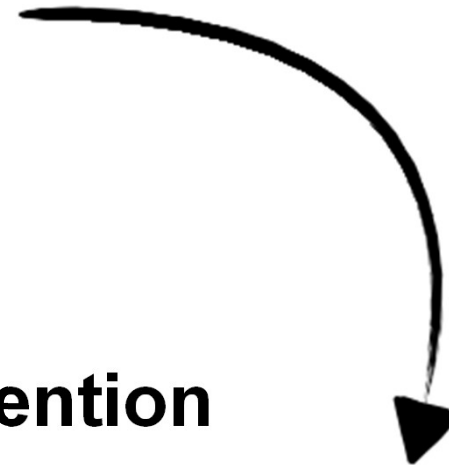
David  
Weiss<sup>2</sup>



Andrew  
McCallum<sup>1</sup>

# Linguistically-Informed Self-Attention

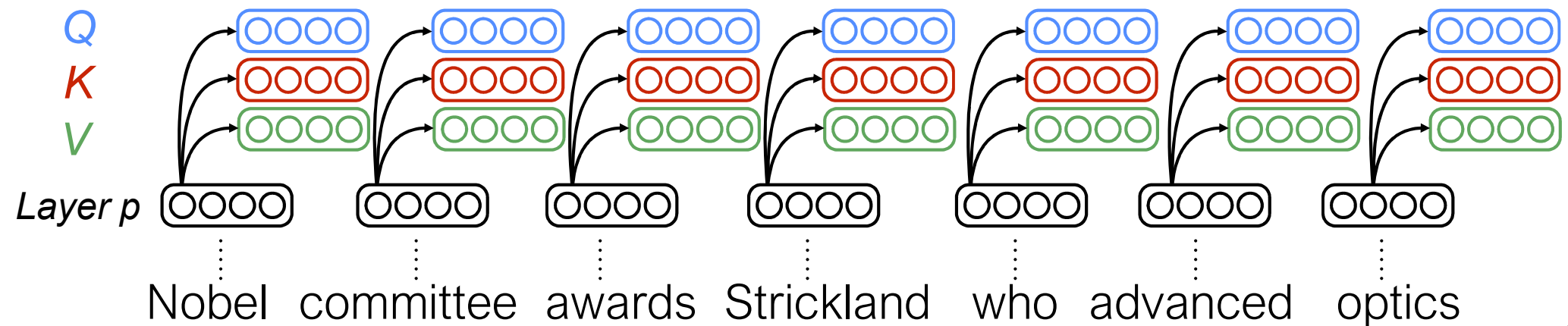
- **Multi-task learning**
  - Part-of-speech tagging
  - Labeled dependency parsing
  - Predicate detection
  - Semantic role spans & labeling
- **Syntactically-informed self-attention**
  - Multi-head self-attention supervised by **syntax**



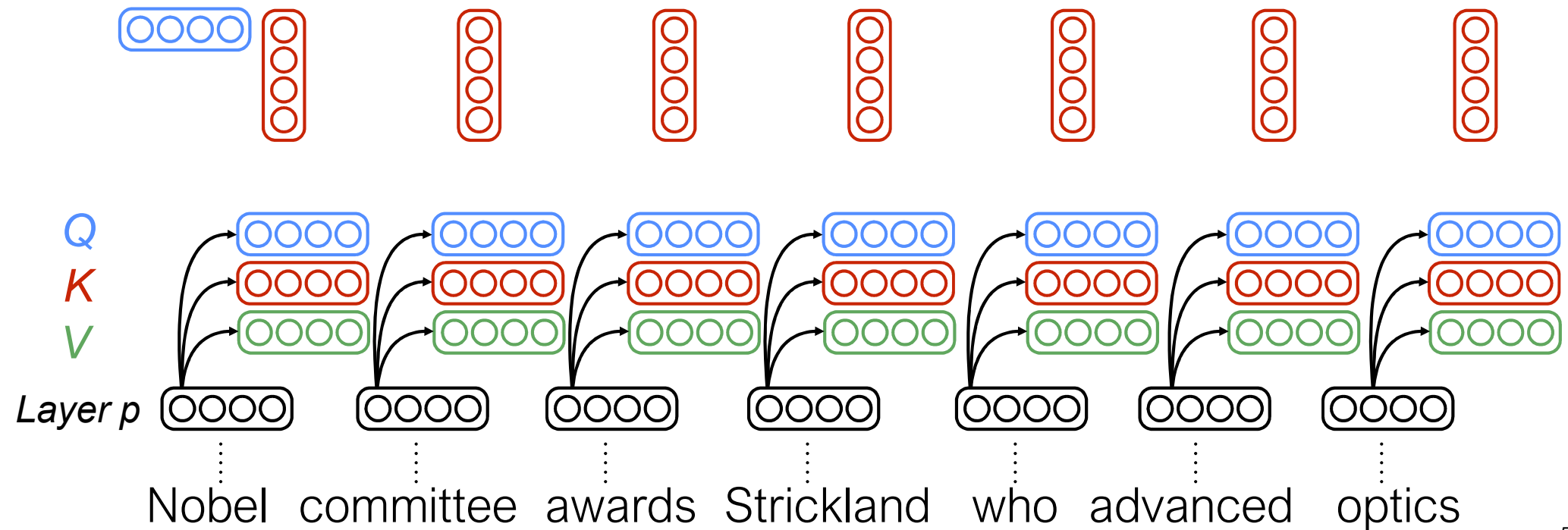
# Outline

- LISA: Linguistically-informed self attention
  - Multi-head self-attention
  - Syntactically-informed self-attention  
[Vaswani et al. 2017]
  - Multi-task learning, single-pass inference
- Experimental results & error analysis

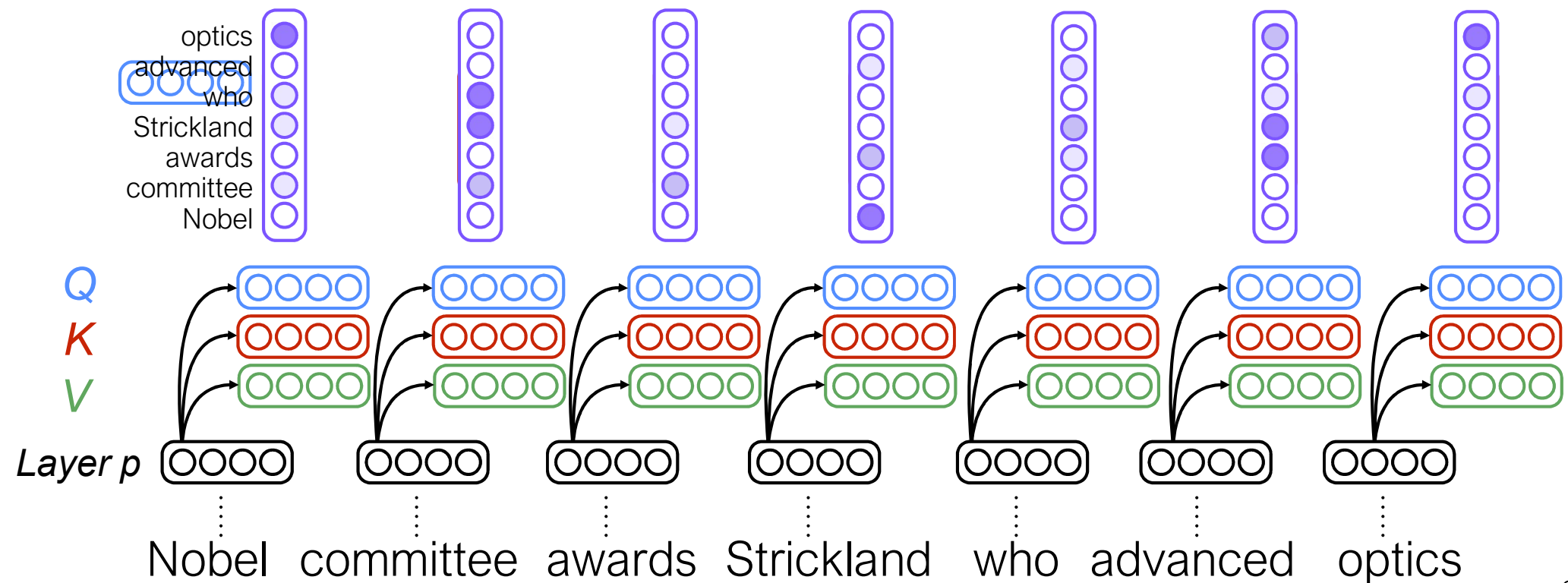
# Self-attention



# Self-attention

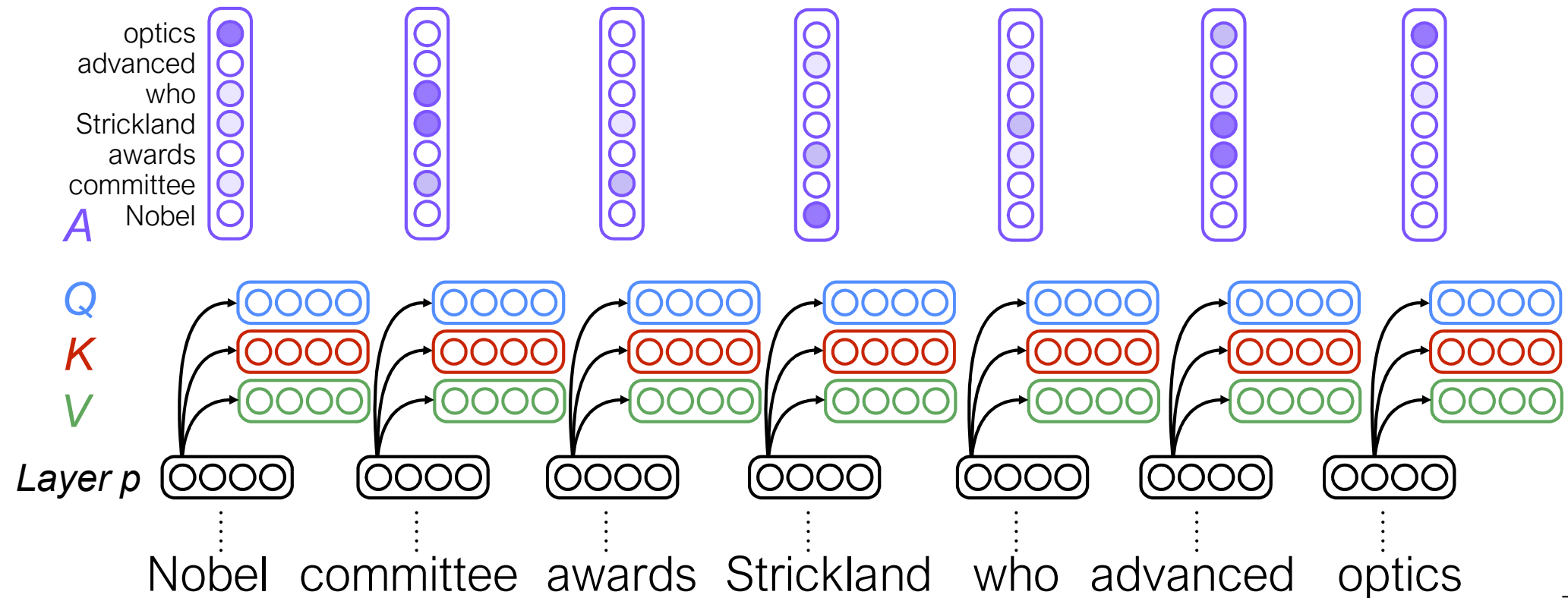


# Self-attention

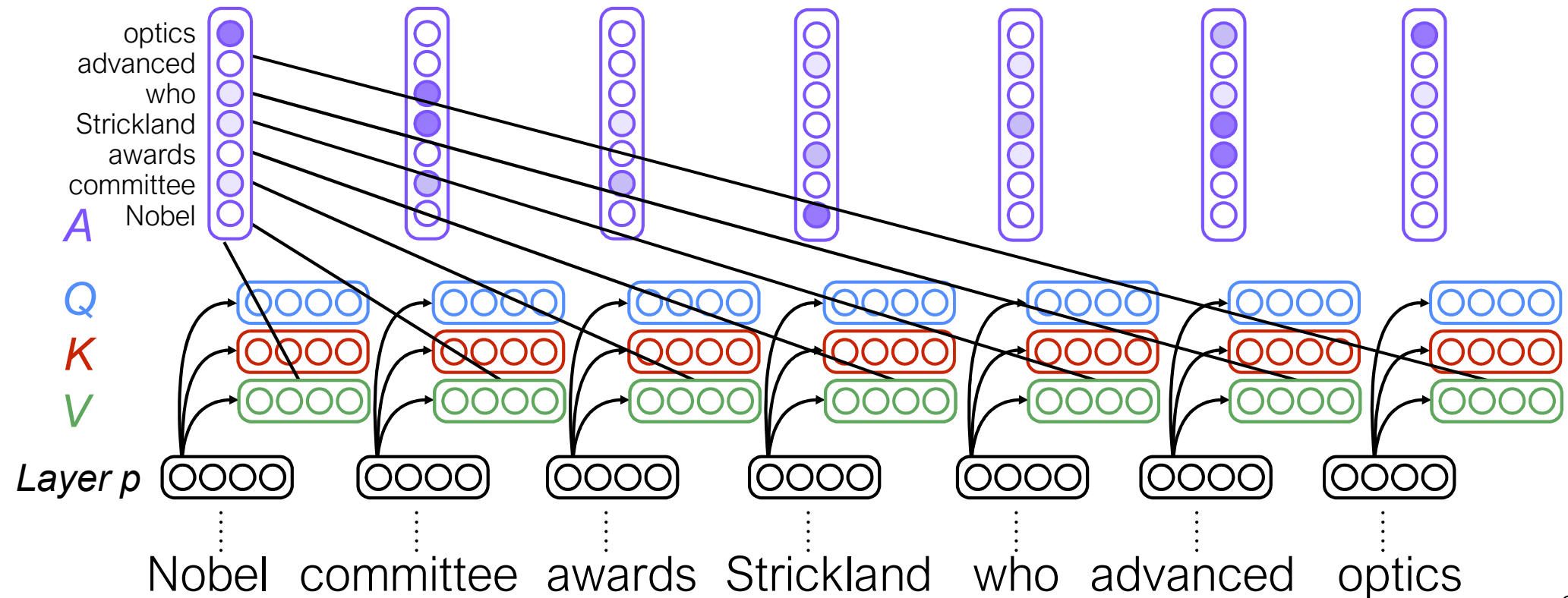




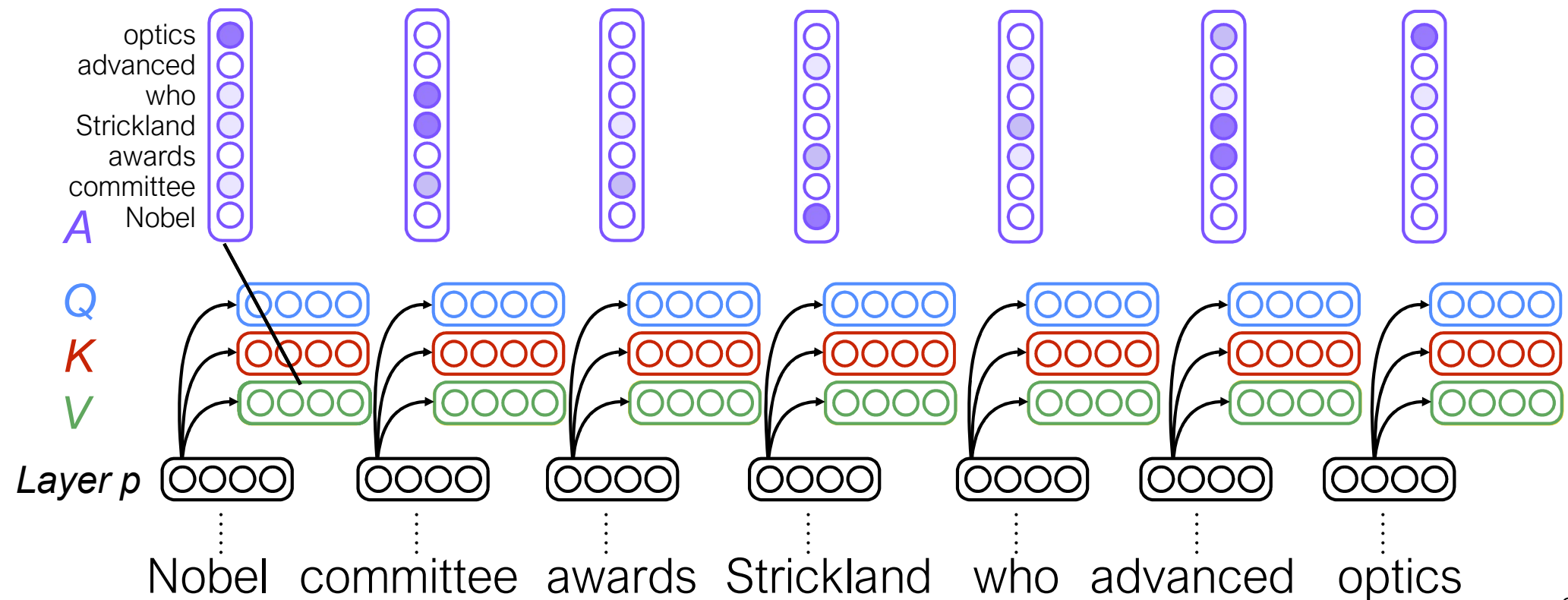
# Self-attention



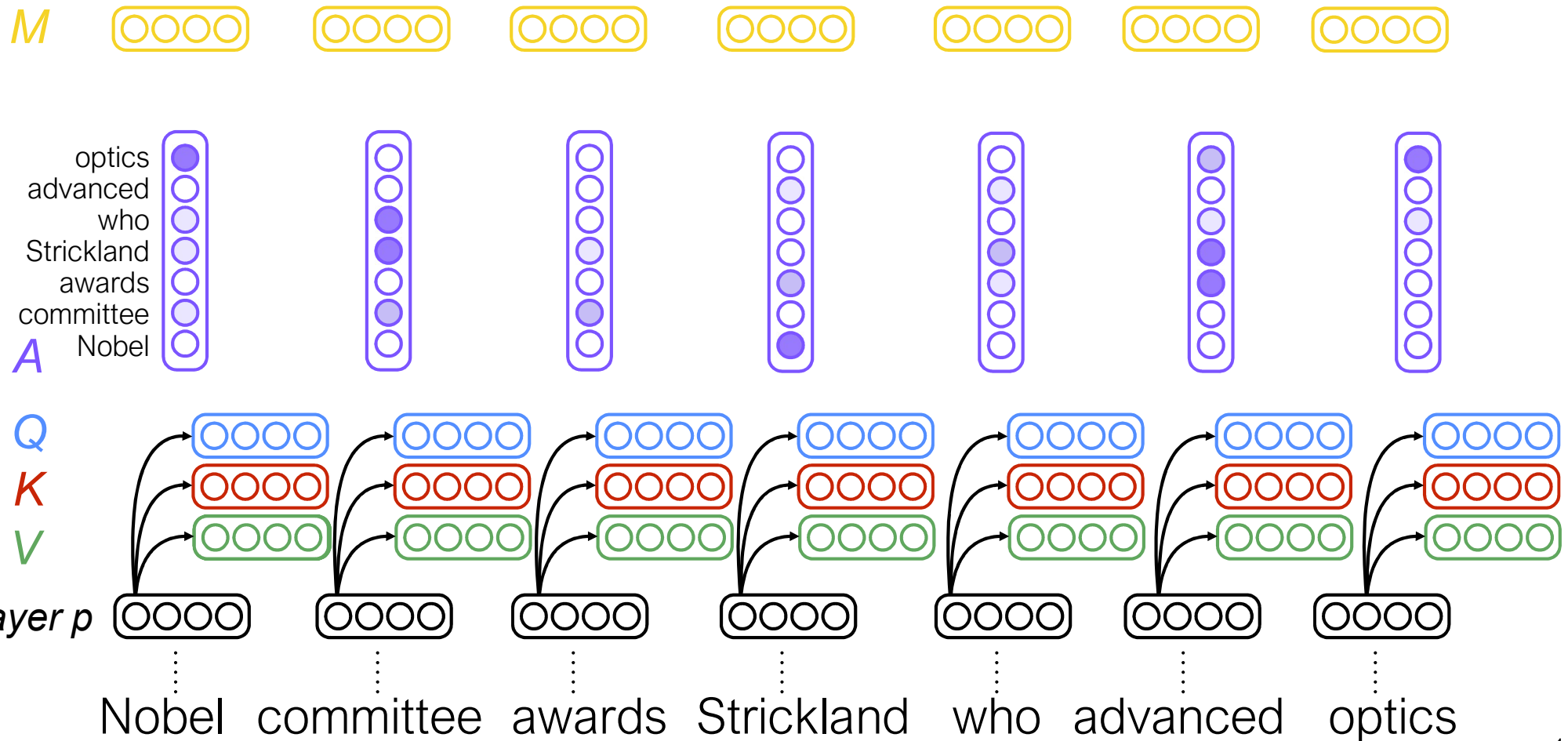
# Self-attention



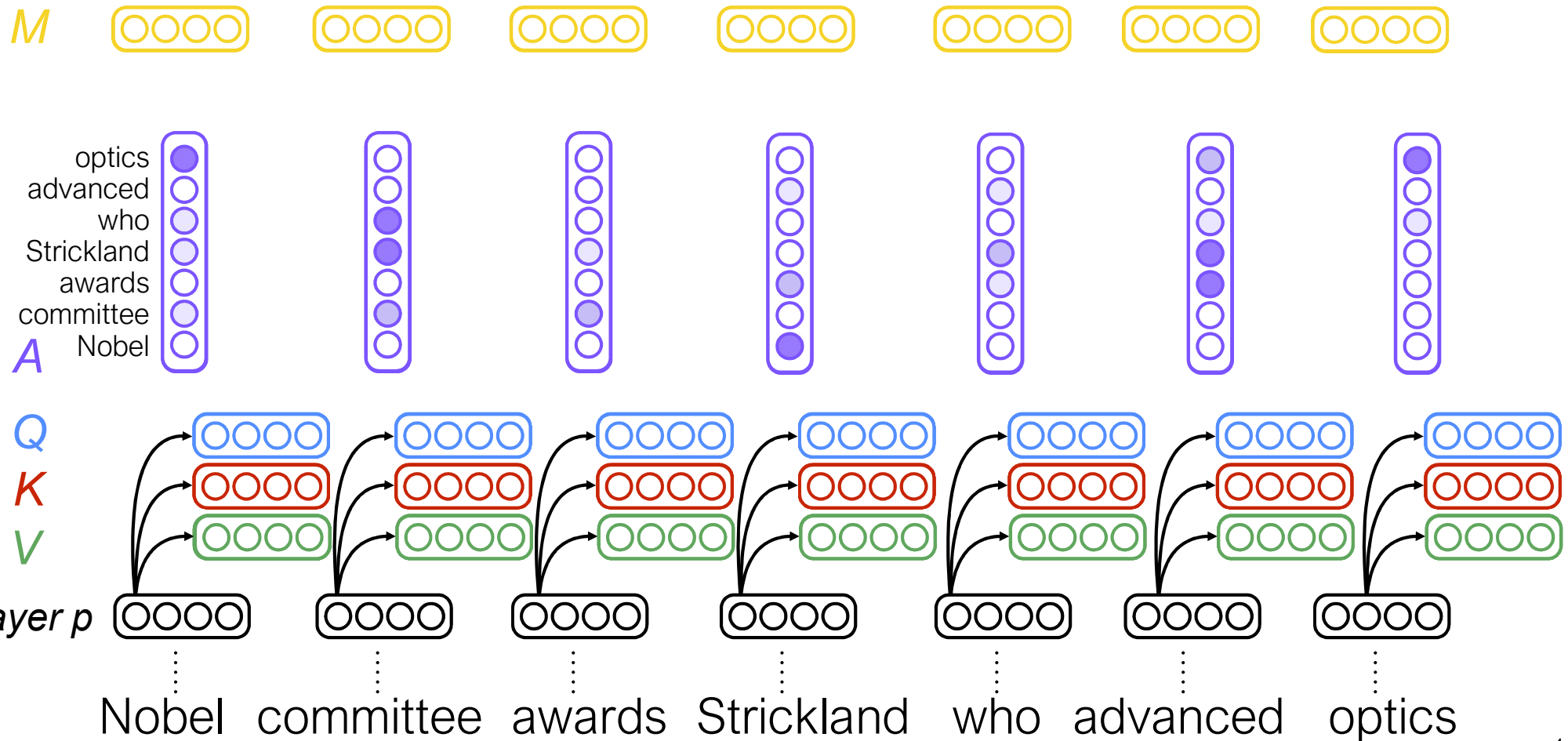
# Self-attention



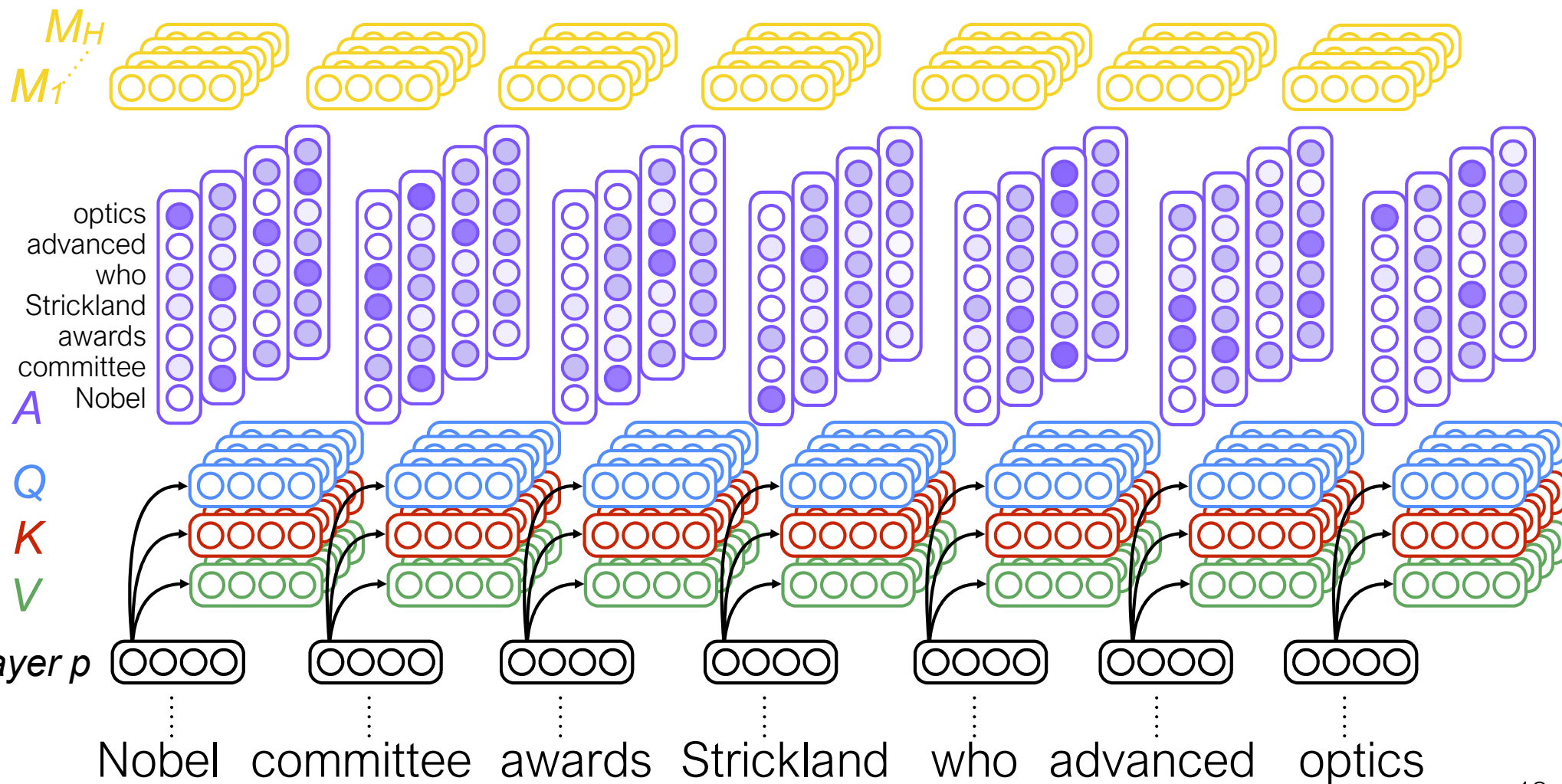
# Self-attention



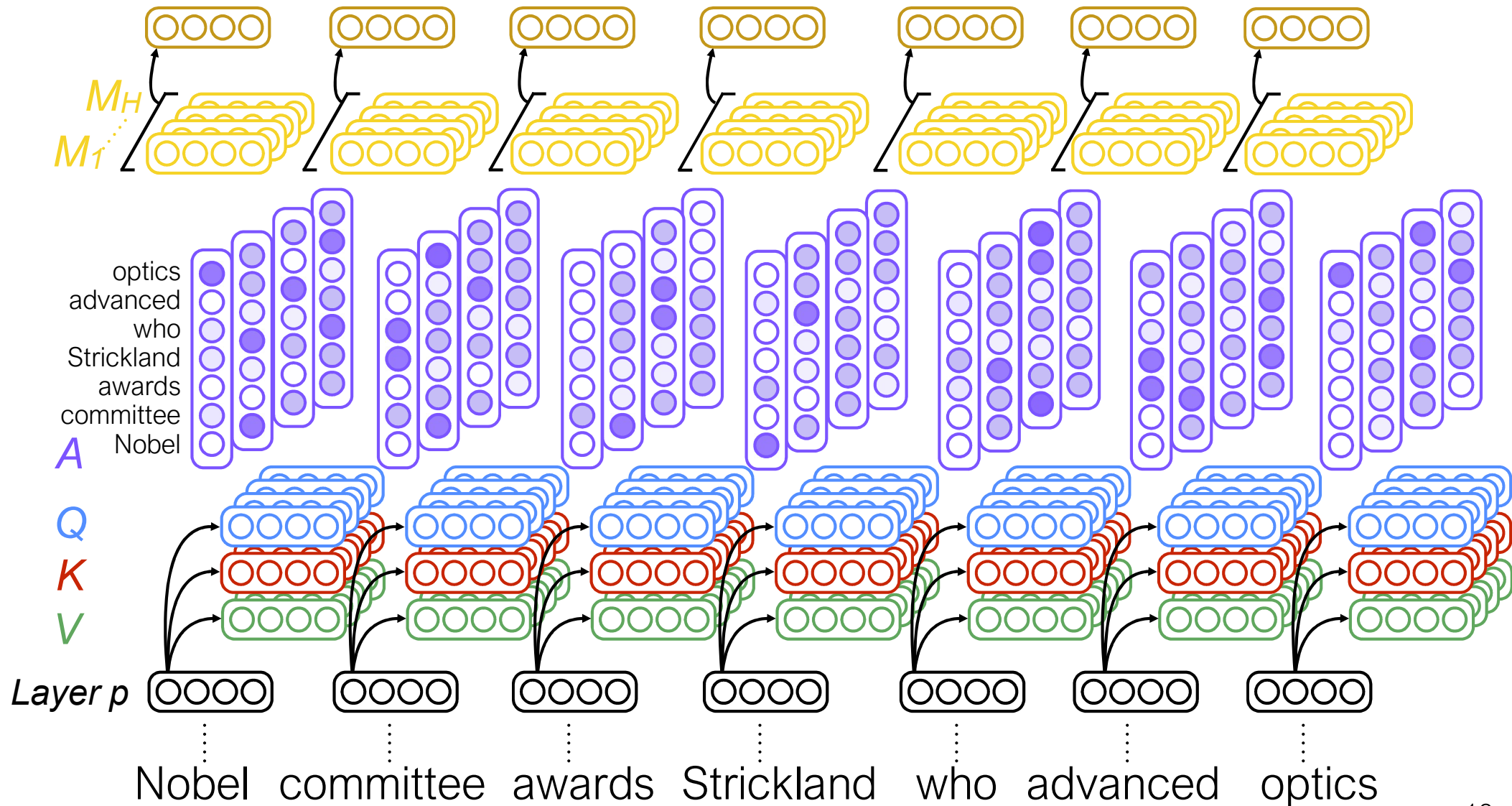
# Self-attention



# Multi-head self-attention

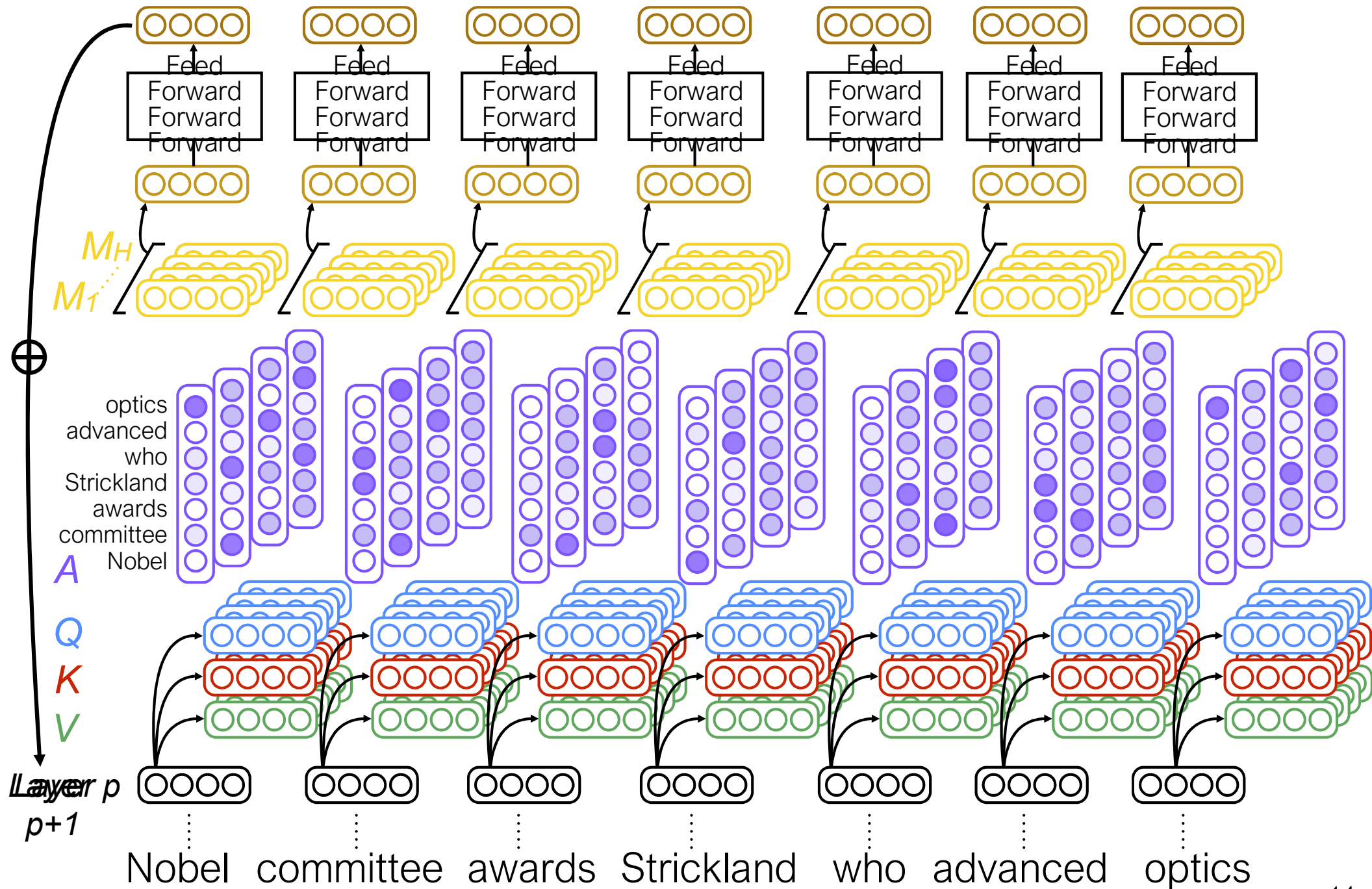


# Multi-head self-attention



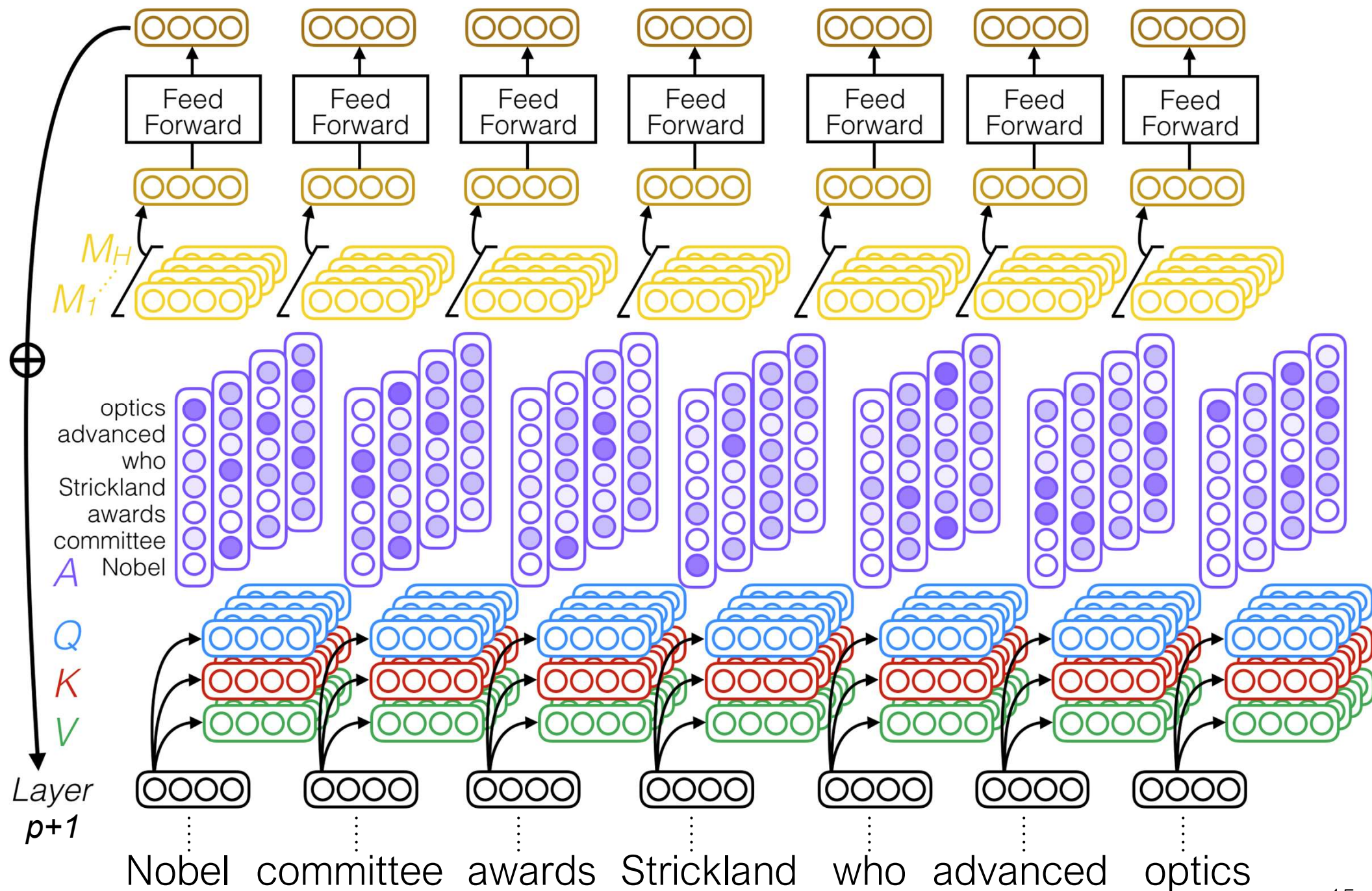


# Multi-head self-attention

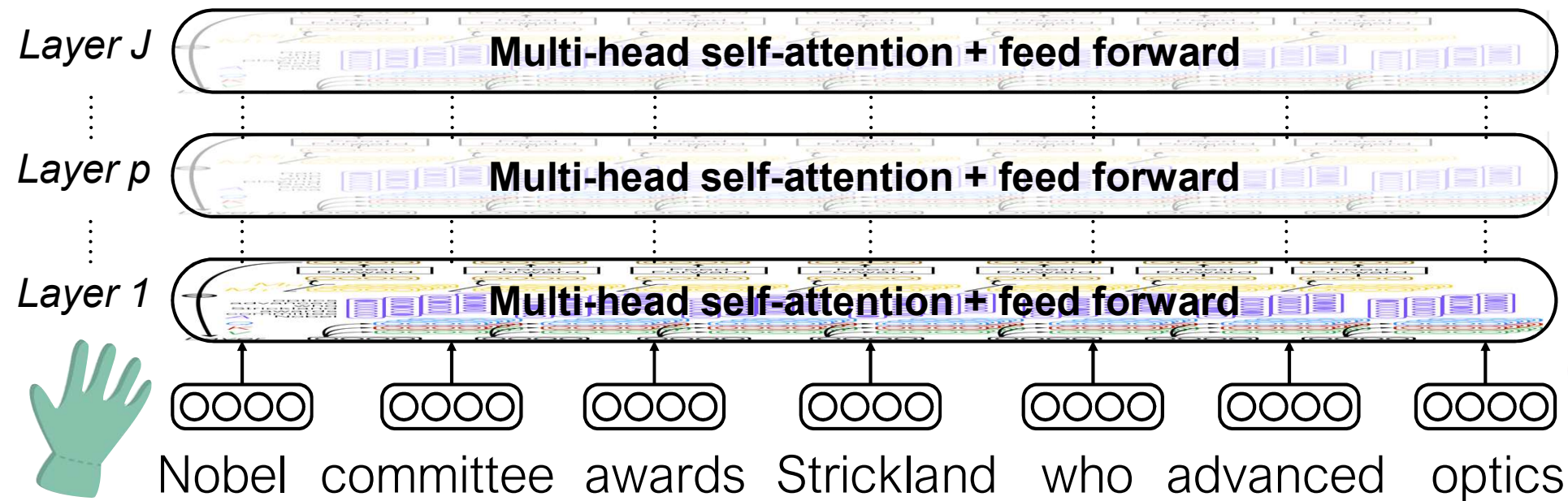




# Multi-head self-attention



# Multi-head self-attention

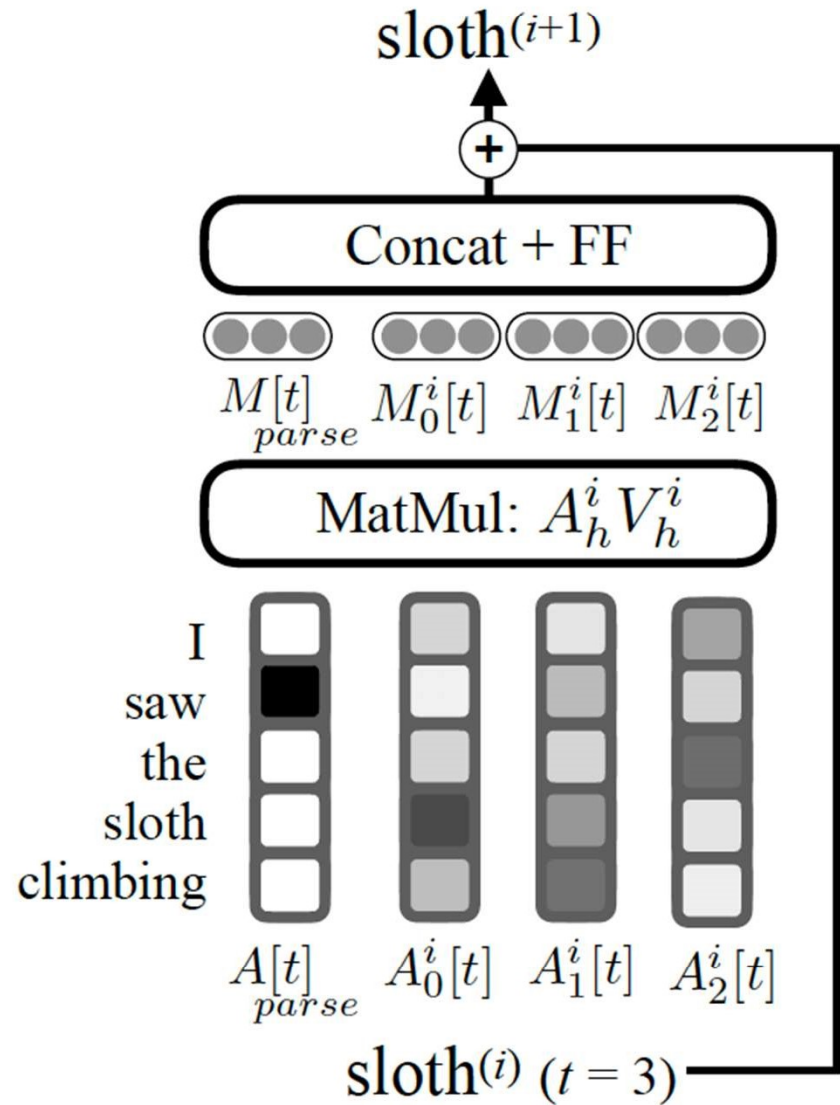


# Self-attention

$$s_t^{(j)} = \text{LN}(s_t^{(j-1)} + T^{(j)}(s_t^{(j-1)}))$$

$$A_h^{(j)} = \text{softmax}(d_k^{-0.5} Q_h^{(j)} K_h^{(j)T})$$

$$M_h^{(j)} = A_h^{(j)} V_h^{(j)}$$



# Outline

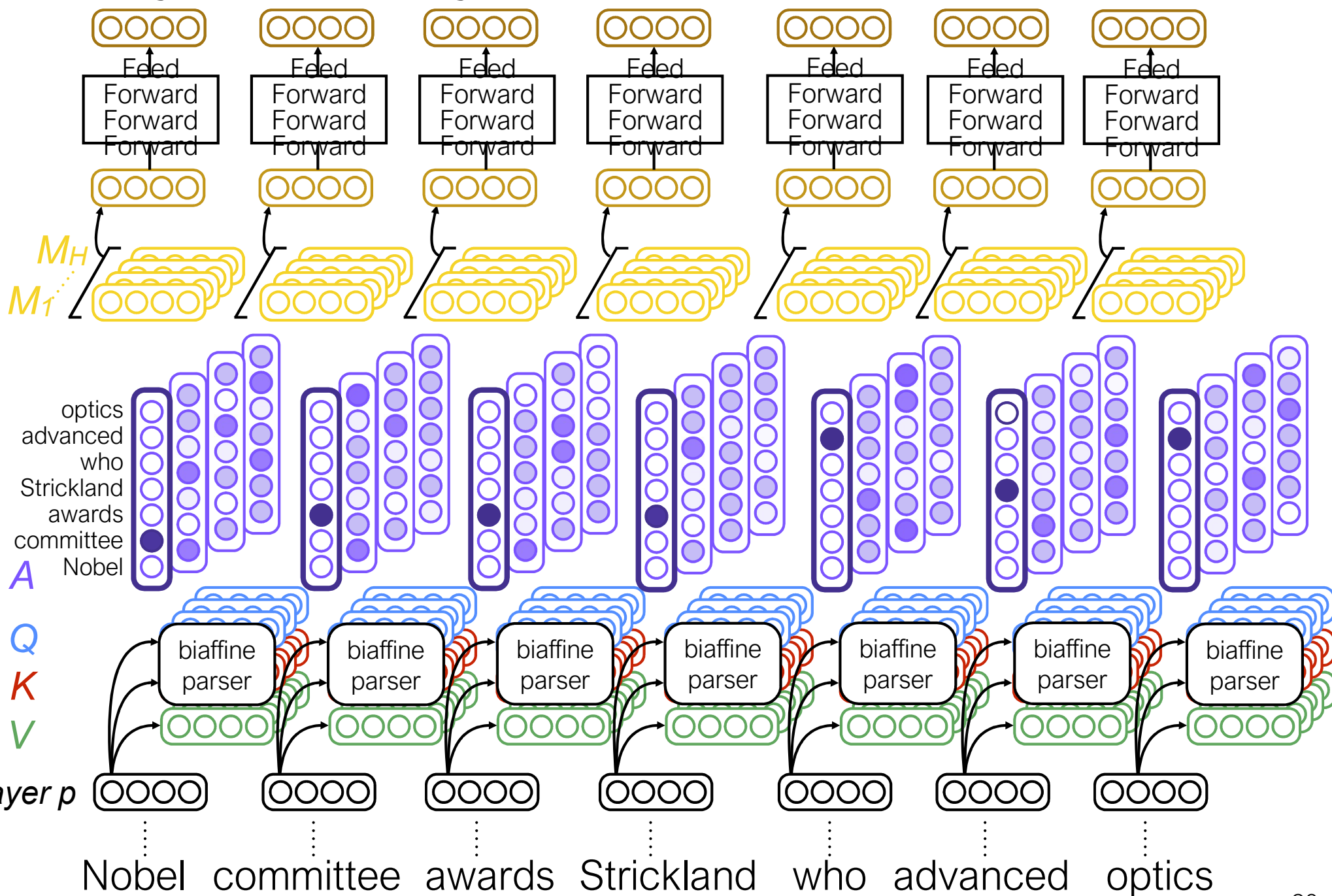
- LISA: Linguistically-informed self attention
  - Multi-head self-attention
  - Syntactically-informed self-attention  
[Vaswani et al. 2017]
  - Multi-task learning, single-pass inference
- Experimental results & error analysis



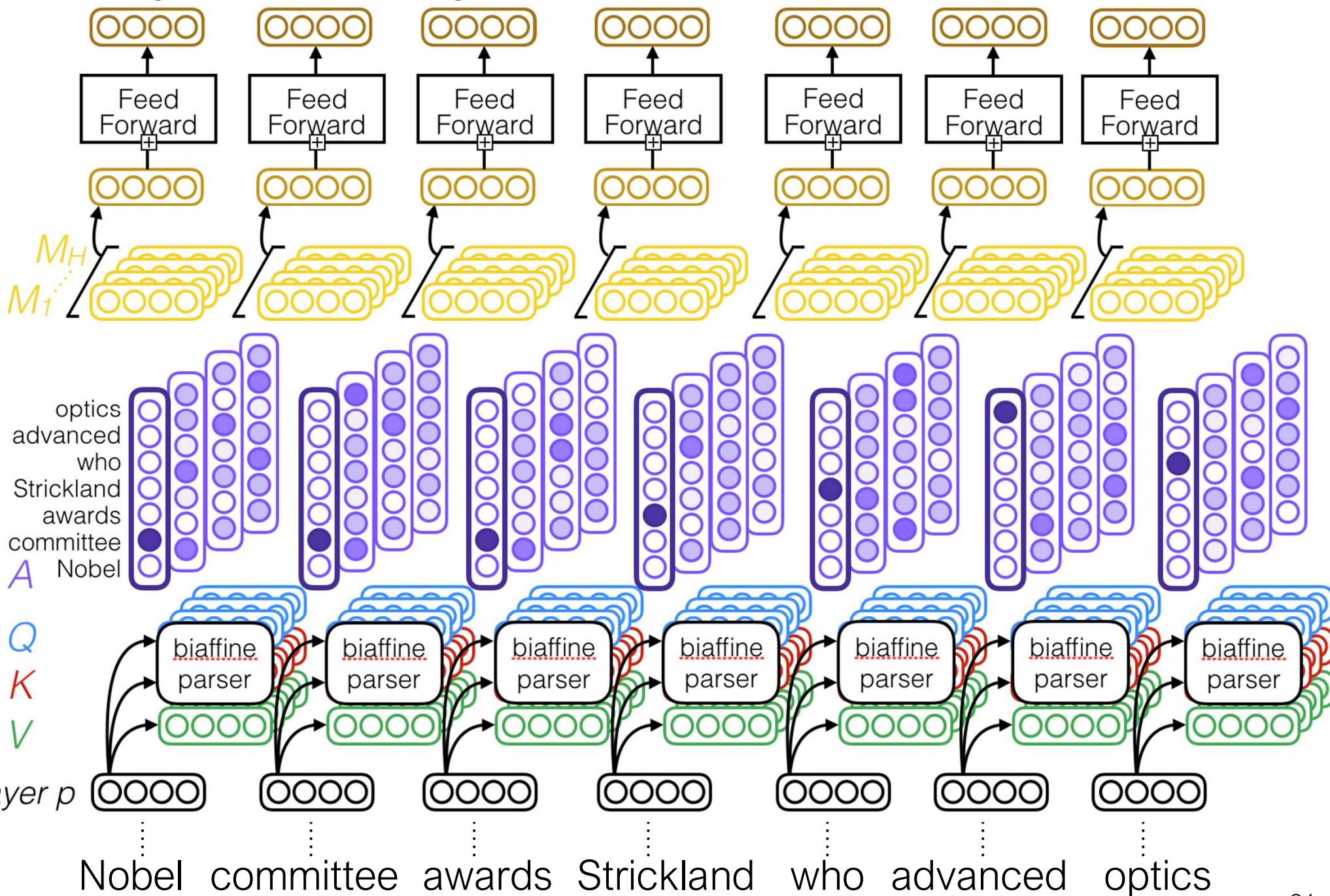
# How to incorporate syntax?

- Multi-task learning [Caruana 1993; Collobert et al. 2011]:
  - Overfits to training domain like single-task end-to-end NN.
  - Must re-train SRL model to leverage new (improved) syntax.
- Dependency path embeddings [Roth & Lapata 2016]; Graph CNN over parse [Marcheggiani & Titov 2017]
  - Restricted context: path to predicate or fixed-width window.
- Syntactically-informed self-attention
  - In one head, token attends to its likely syntactic parent(s).
  - Global context: In next layer, tokens observe all other parents.
  - At test time: can use own predicted parse, **OR** supply syntax to improve SRL model without re-training.

# Syntactically-informed self-attention

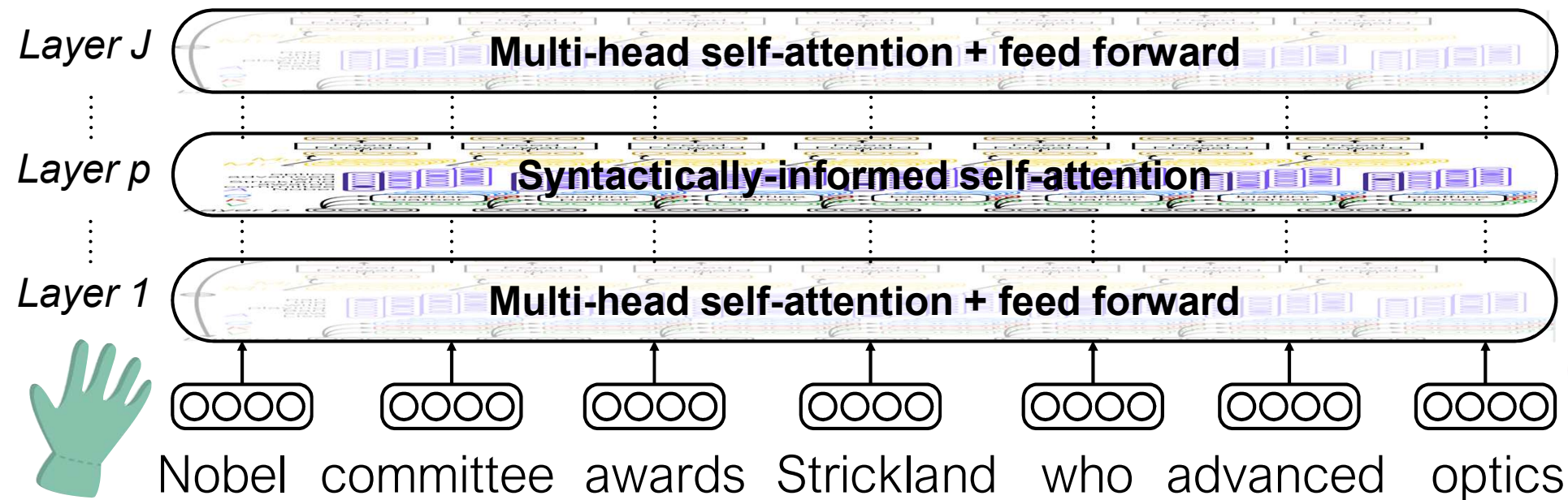


# Syntactically-informed self-attention





# Syntactically-informed self-attention

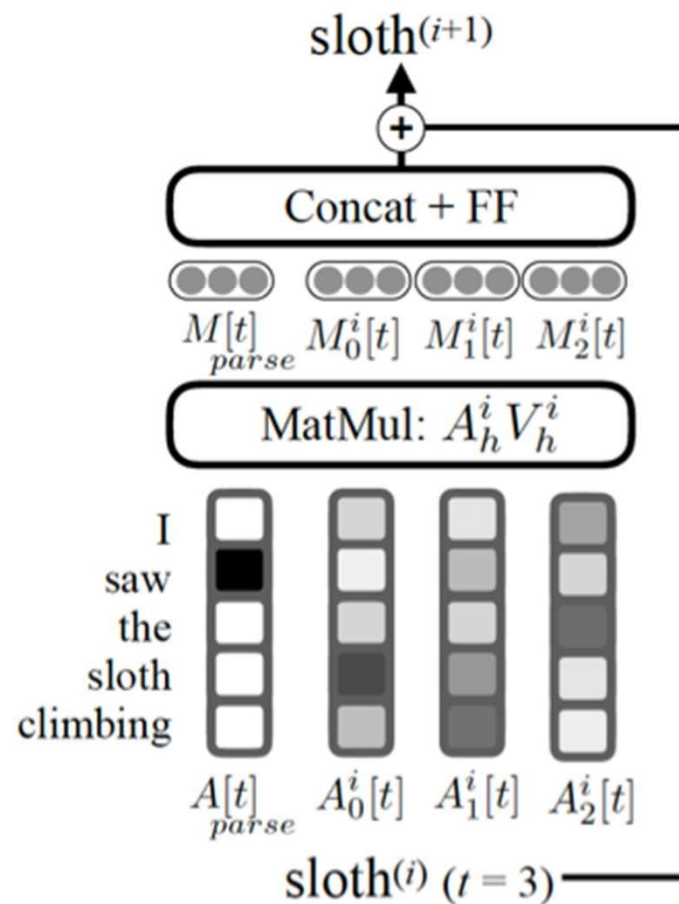




# Syntactically-informed self-attention

$$A_{parse} = \text{softmax}(Q_{parse} U_{heads} K_{parse}^T)$$

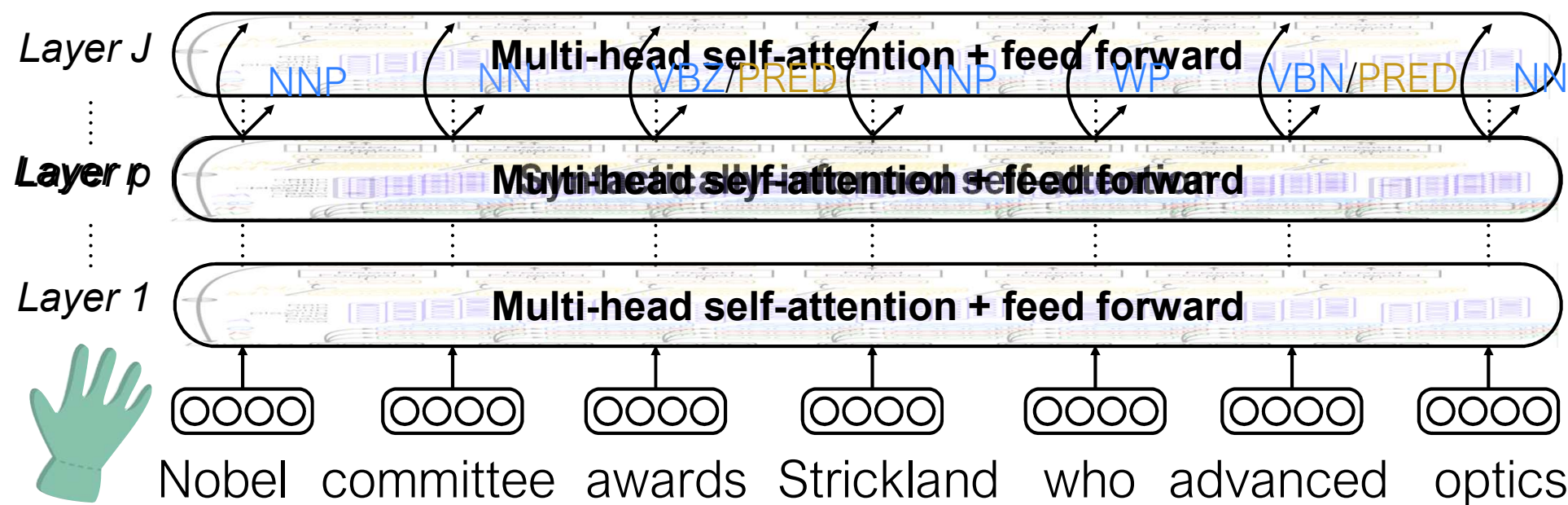
$$P(q = \text{head}(t) \mid \mathcal{X}) = A_{parse}[t, q]$$



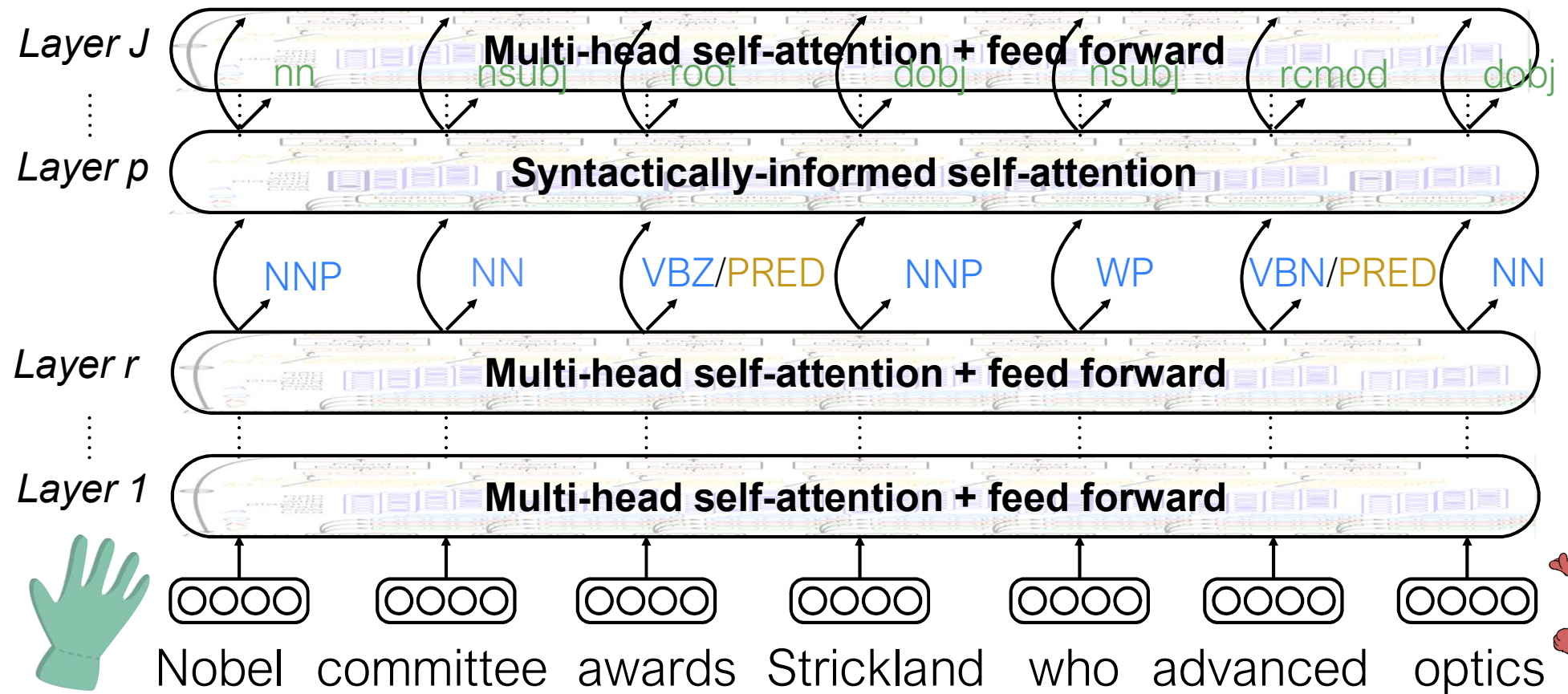
# Outline

- LISA: Linguistically-informed self attention
  - Multi-head self-attention
  - Syntactically-informed self-attention
  - Multi-task learning, single-pass inference
- Experimental results & error analysis

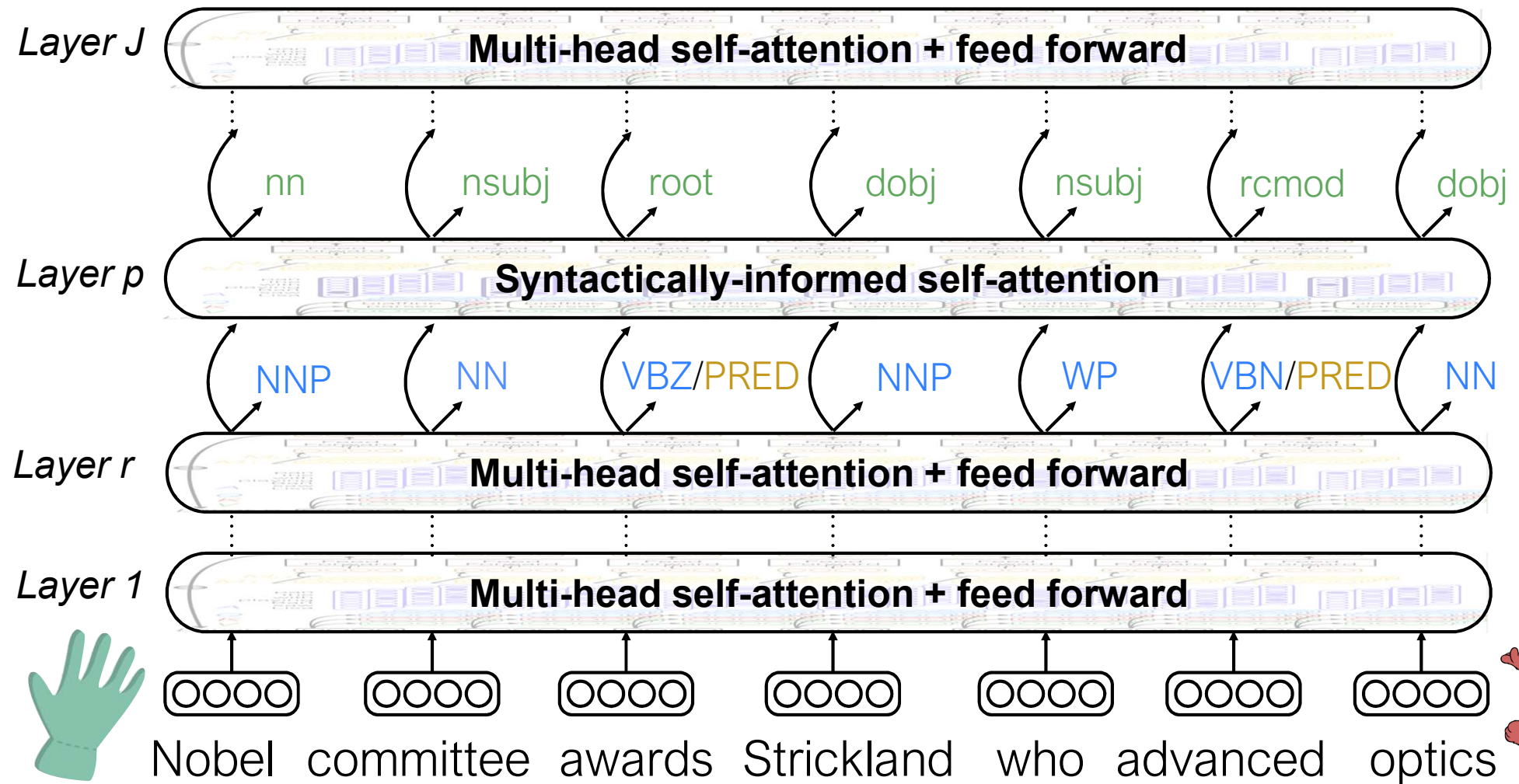
# LISA: Linguistically-Informed Self-Attention



# LISA: Linguistically-Informed Self-Attention

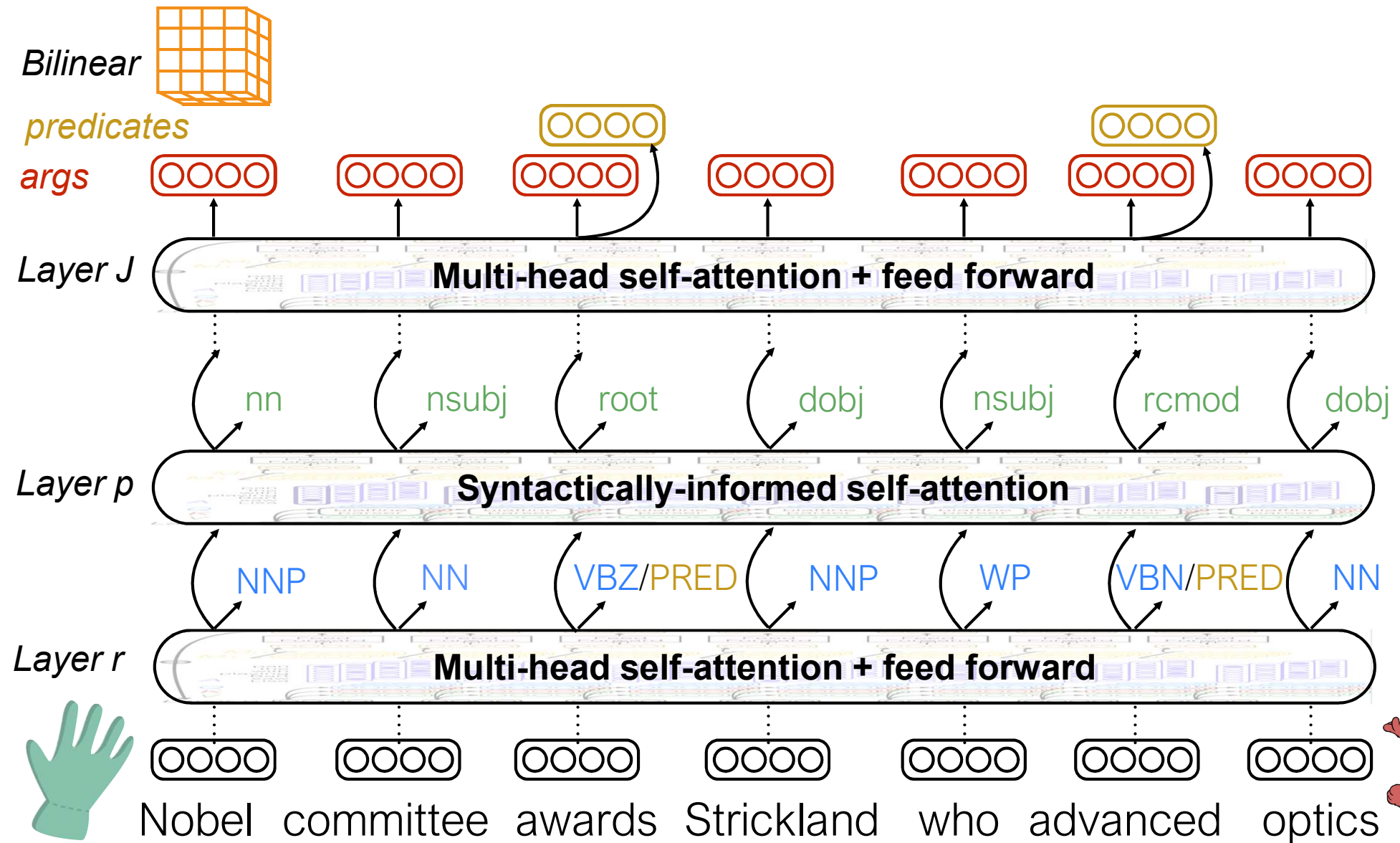


# LISA: Linguistically-Informed Self-Attention

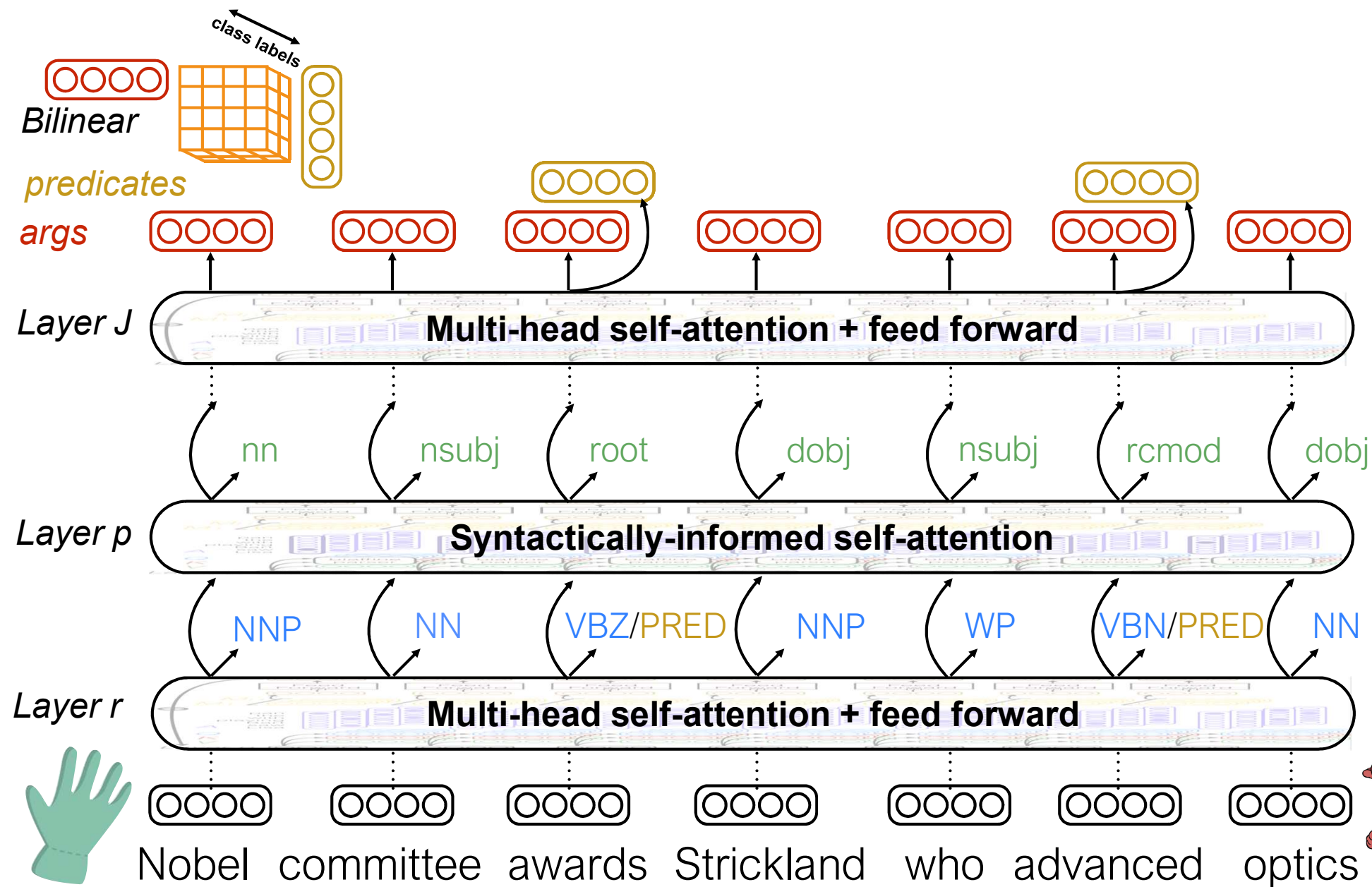




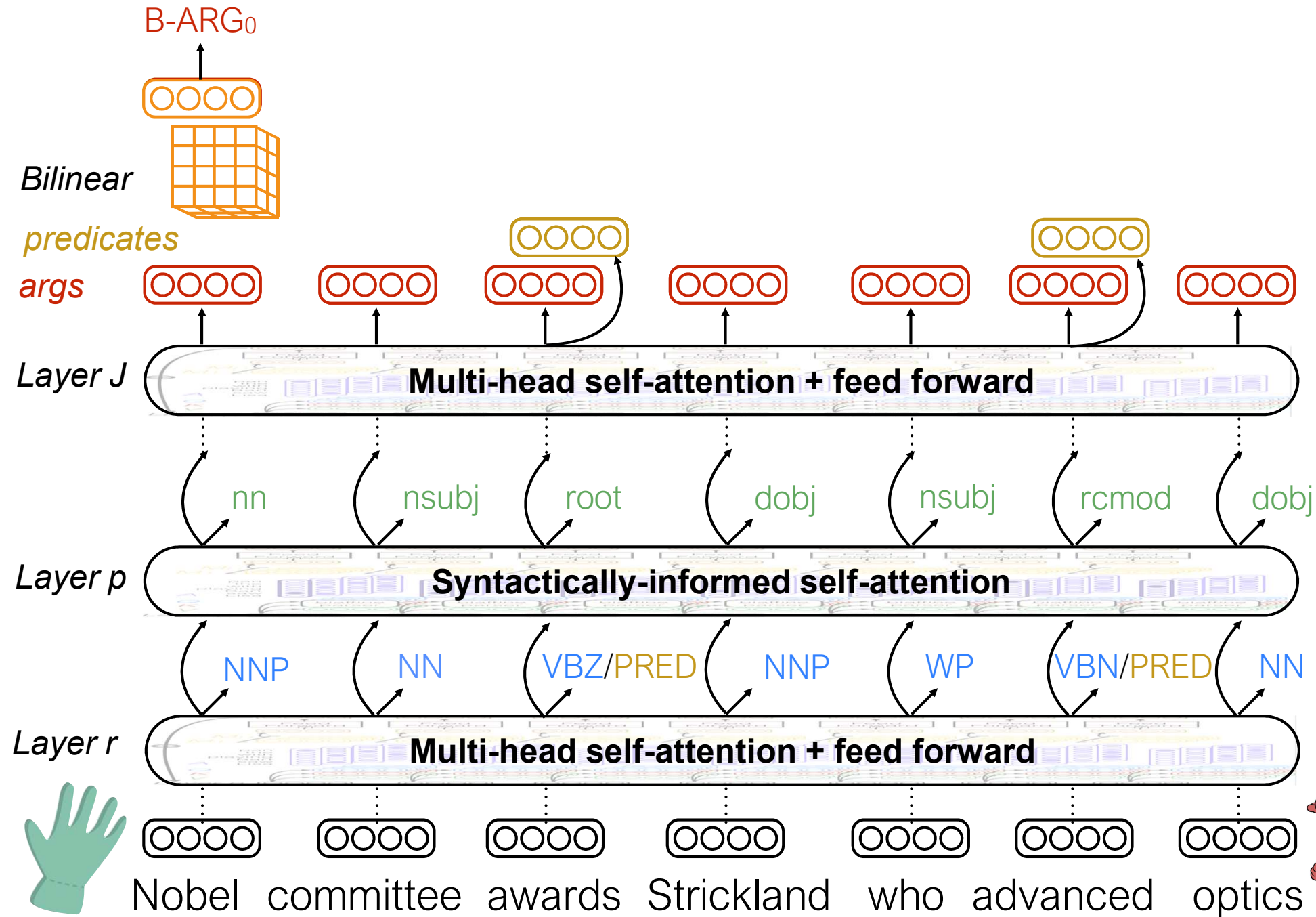
# LISA: Linguistically-Informed Self-Attention



# LISA: Linguistically-Informed Self-Attention

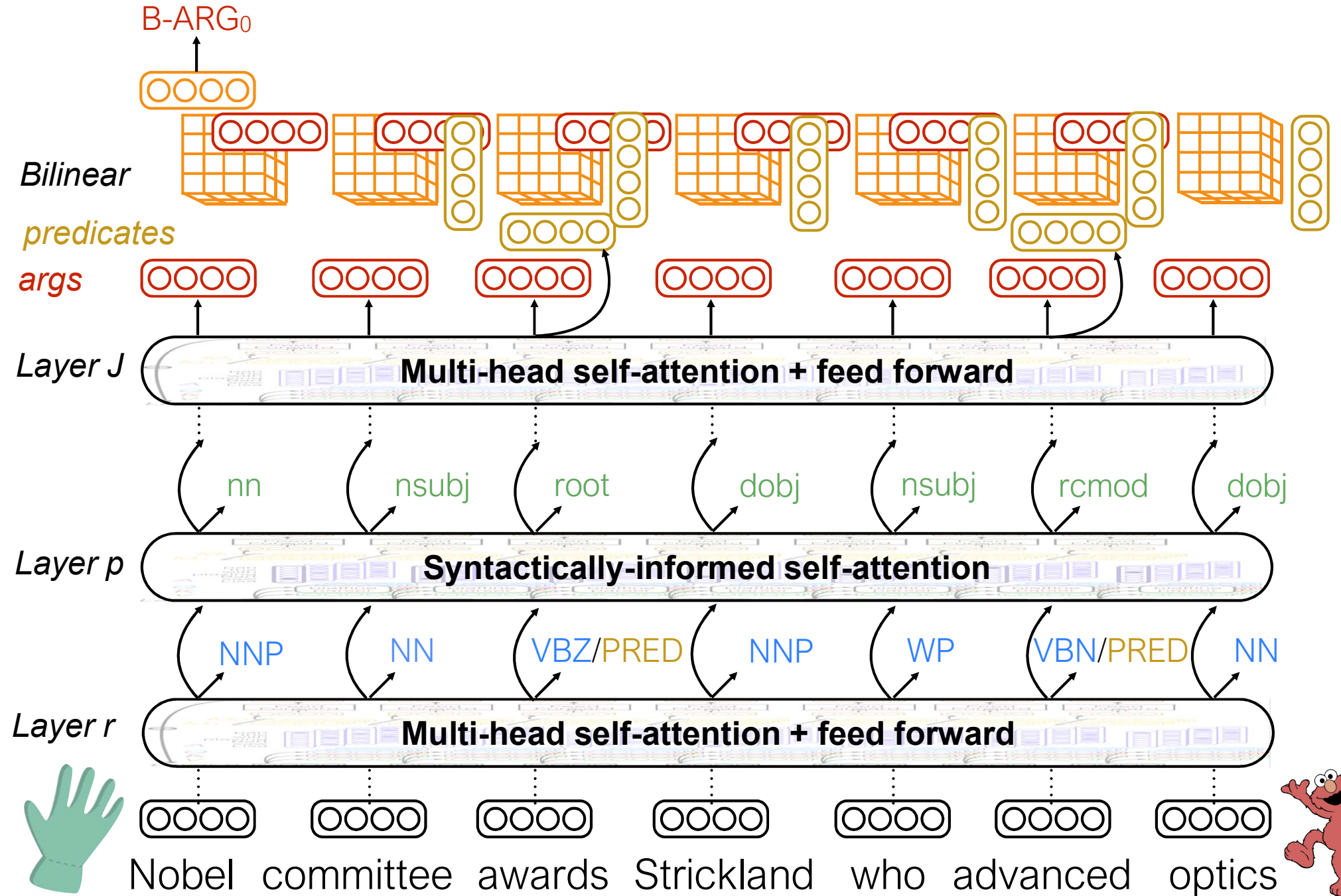


# LISA: Linguistically-Informed Self-Attention

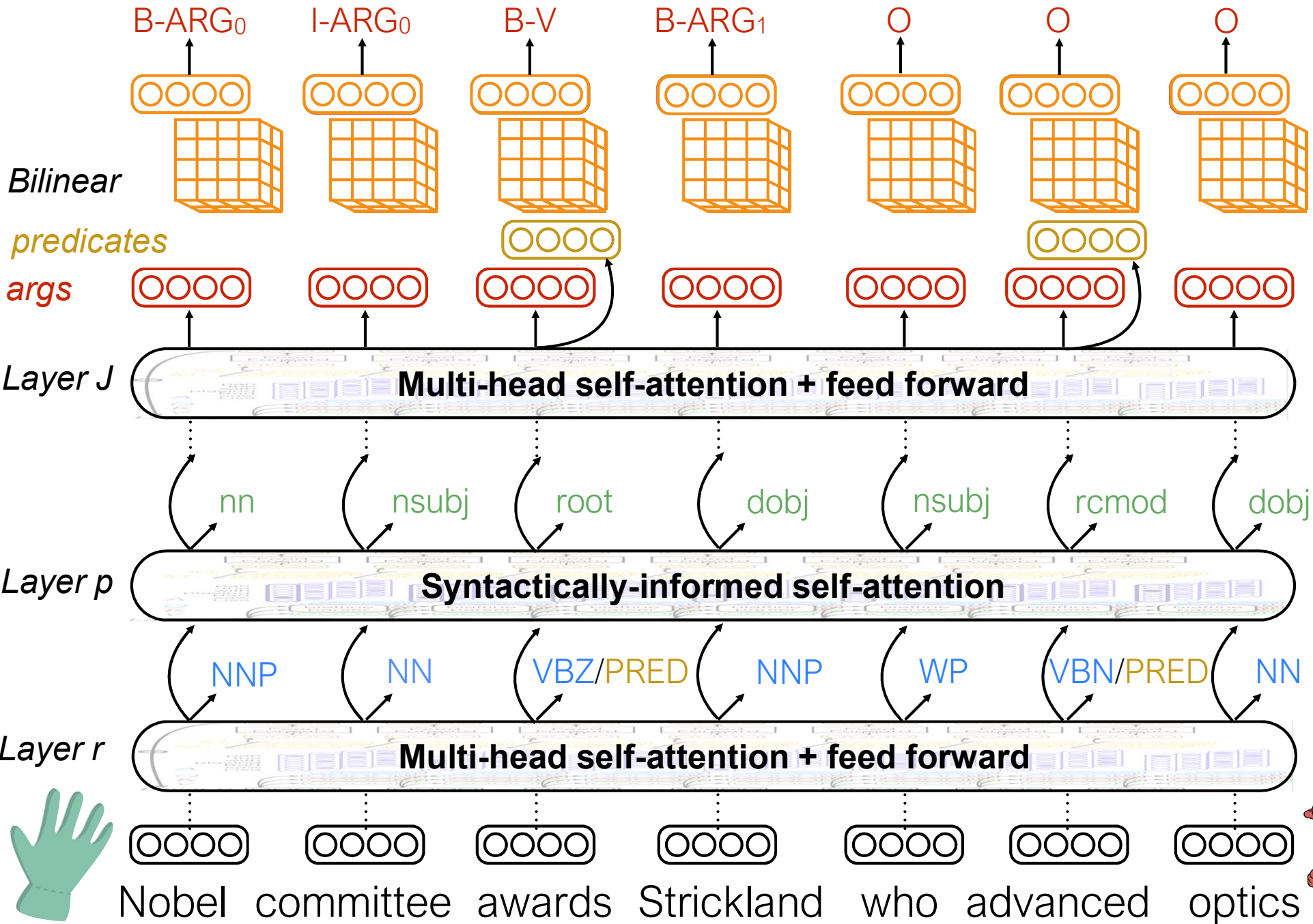




# LISA: Linguistically-Informed Self-Attention

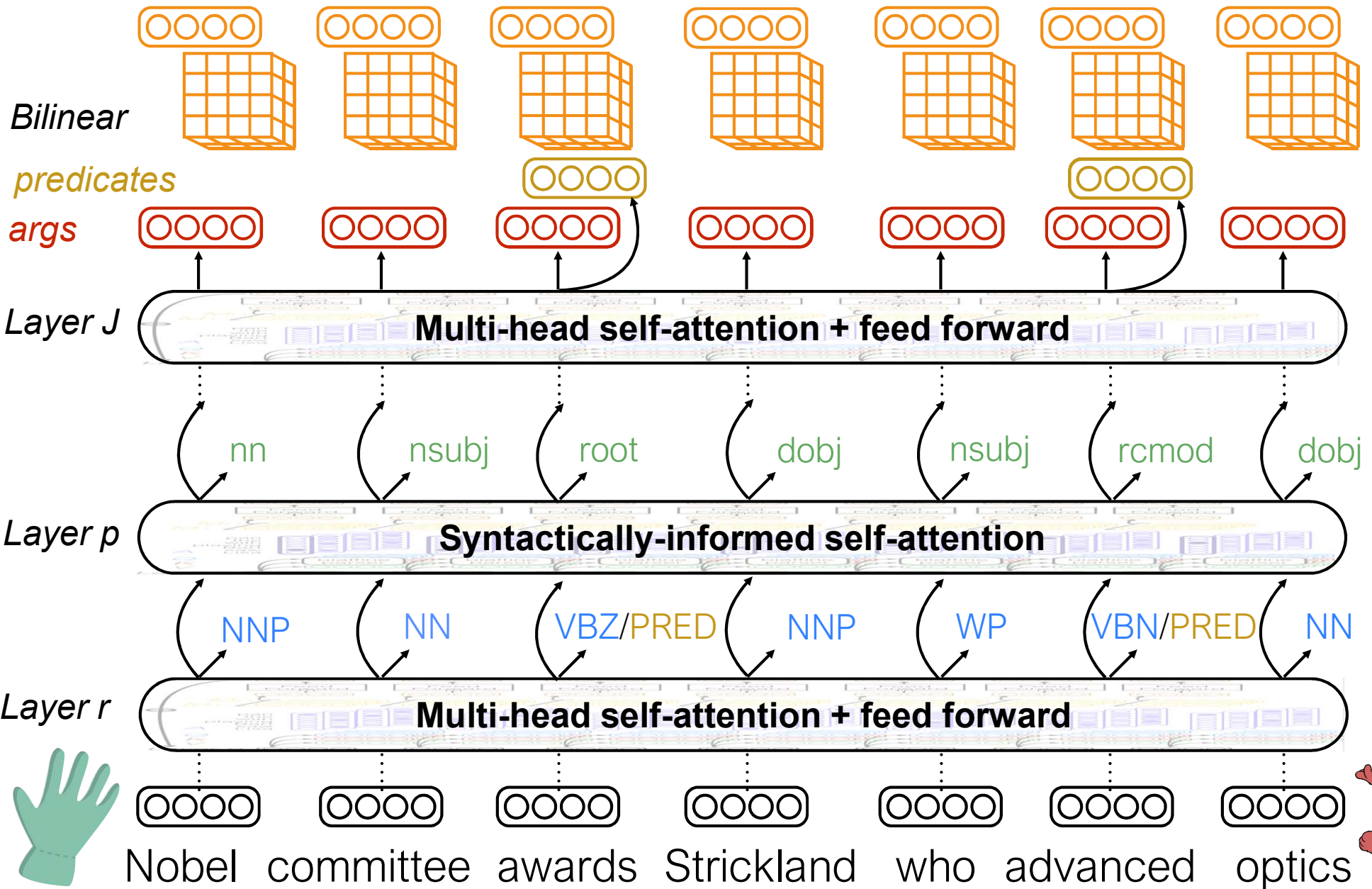


# LISA: Linguistically-Informed Self-Attention



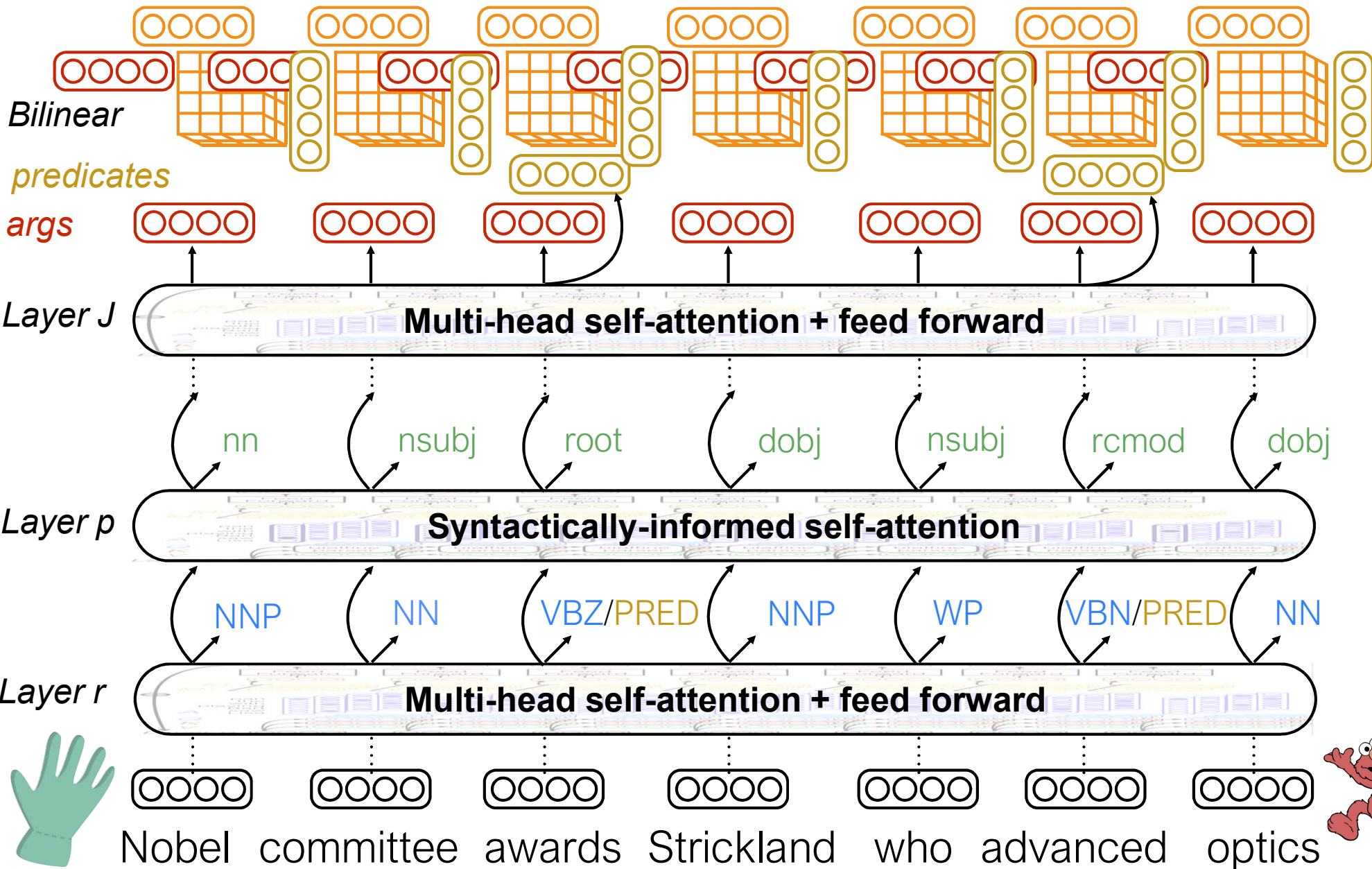
# LISA: Linguistically-Informed Self-Attention

B-ARG<sub>0</sub> I-ARG<sub>0</sub> B-V B-ARG<sub>1</sub> O O O



# LISA: Linguistically-Informed Self-Attention

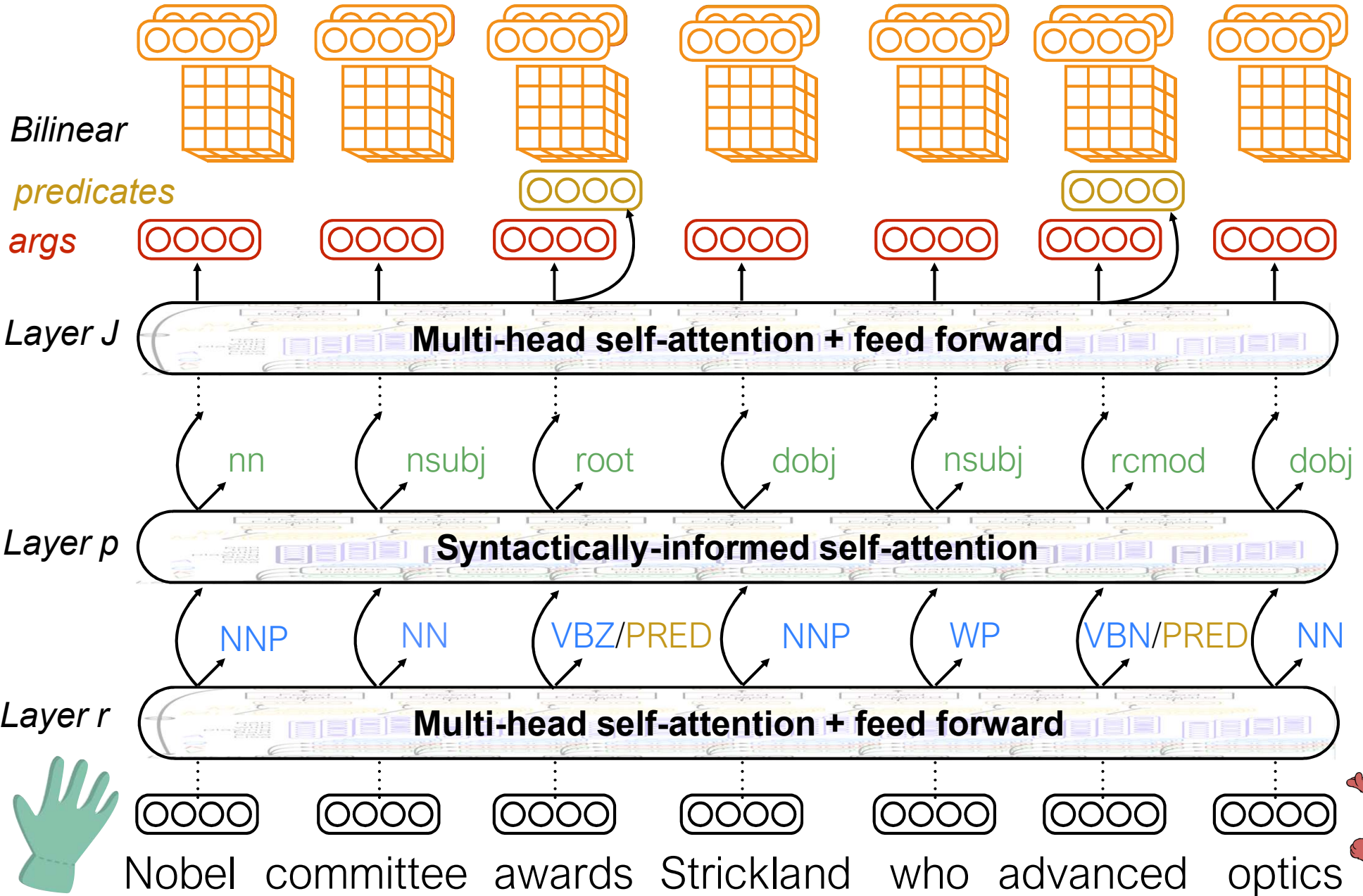
B-ARG<sub>0</sub> I-ARG<sub>0</sub> B-V B-ARG<sub>1</sub> O O O



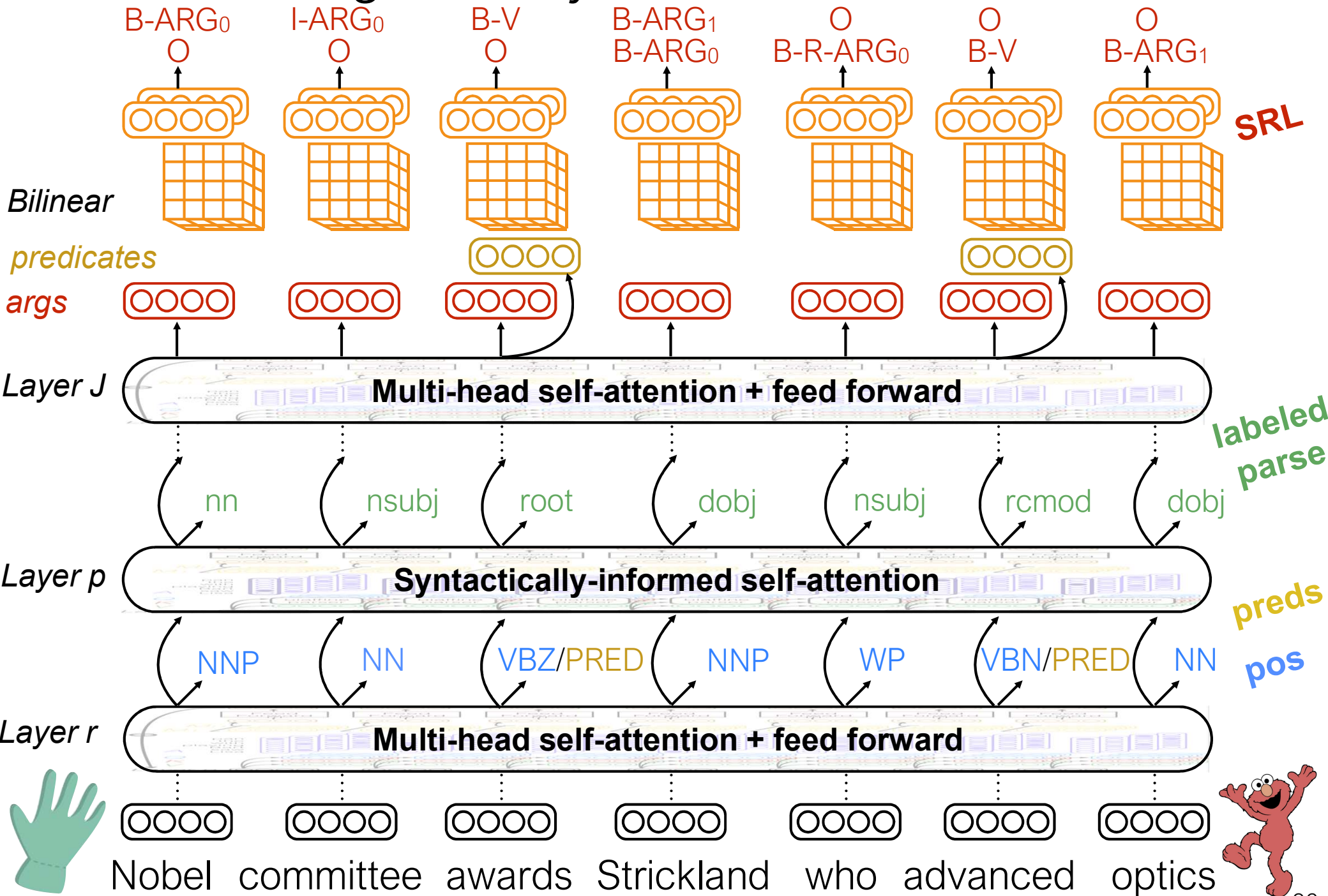


# LISA: Linguistically-Informed Self-Attention

B-ARG<sub>0</sub> I-ARG<sub>0</sub> B-V B-ARG<sub>1</sub> O O O



# LISA: Linguistically-Informed Self-Attention



# LISA: Linguistically-Informed Self-Attention

$$\frac{1}{T} \sum_{t=1}^T \left[ \sum_{f=1}^F \log P(y_{ft}^{role} \mid \mathcal{P}_G, \mathcal{V}_G, \mathcal{X}) \right. \\ \left. + \log P(y_t^{prp} \mid \mathcal{X}) \right. \\ \left. + \lambda_1 \log P(\text{head}(t) \mid \mathcal{X}) \right. \\ \left. + \lambda_2 \log P(y_t^{dep} \mid \mathcal{P}_G, \mathcal{X}) \right]$$

# Outline

- LISA: Linguistically-informed self attention
  - Multi-head self-attention
  - Syntactically-informed self-attention
  - Multi-task learning, single-pass inference
- Experimental results & error analysis





# Experimental results

	CoNLL-2005	CoNLL-2012
<b>domains</b>	Train, dev: news Test: news, novels	Train, dev, test: 7 domains (news, telephone, bible, ...)
<b>word embeddings</b>	GloVe [Pennington et al. 2014] ELMo [Peters et al. 2018]	GloVe [Pennington et al. 2014] ELMo [Peters et al. 2018]
<b>predicates</b>	predicted; gold	predicted
<b>baselines</b>	He et al. 2017 He et al. 2018 Tan et al. 2018	He et al. 2018
<b>our models</b>	SA LISA LISA+D&M, +Gold <sub>[SEP]</sub> Lisa_Gold	SA LISA LISA+D&M, +Gold <sub>[SEP]</sub> Lisa_Gold

# Experimental results

He et al. 2017	PoE
He et al. 2018	jointly predict all predicates and argument spans
SA	does not incorporate syntactic information
LISA	Predicted parser
+D&M	injecting state-of-the-art predicted parses at test time (+D&M)
<i>+Gold</i>	the gold syntactic parse at test time (+Gold)

# Experimental results

	 <b>GloVe</b>		 <b>ELMo</b>	
	<b>in-domain</b>	<b>out-of-domain</b>	<b>in-domain</b>	<b>out-of-domain</b>
He et al. 2017	82.7	70.1	---	---
He et al. 2018	82.5	70.8	86.0	76.1
SA	83.72	71.51	86.09	76.35
LISA	83.61	71.91	86.55	78.05
+D&M	84.99 <b>94.9 UAS</b>	74.66 <b>90.3 UAS</b>	86.90 <b>96.3 UAS</b>	78.25 <b>96.4 UAS</b>
	+2.49 F1	+3.86 F1	+0.9 F1 ?	+2.15 F1

# Experimental results: CoNLL-2005



GloVe

ELMo

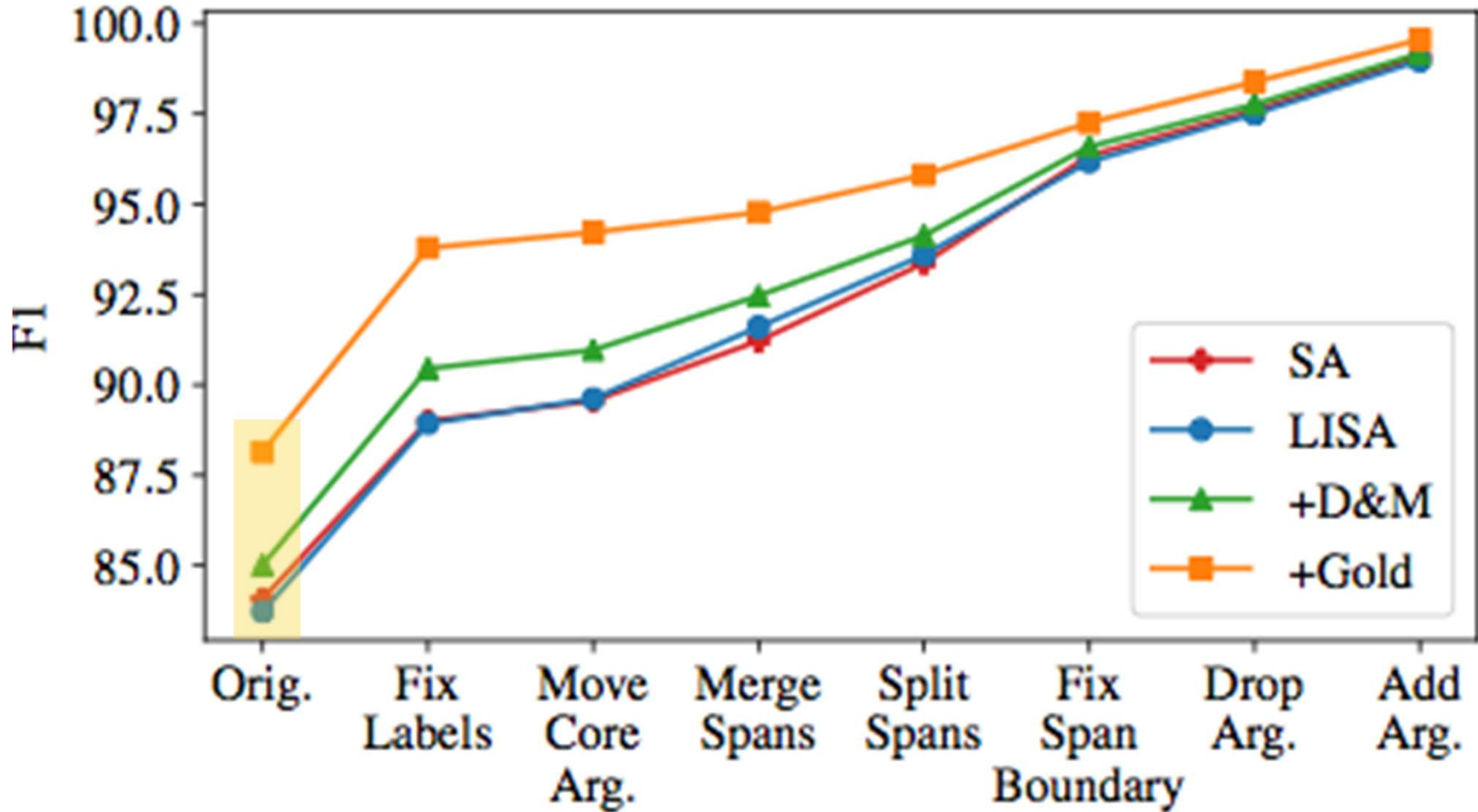


in-domain (dev)

in-domain (dev)

**96.5 UAS!**

# Experimental results: Analysis



# Error Breakdown

Labeling Errors

PP Attachment

Long-range Dependencies

Structural Consistency

Can Syntax Still Help?

## Oracle Transformations

1) Fix Label:  
[We] *fly* to NYC tomorrow.  
~~ARG0~~  
ARG1

2) Move core arg:  
ARG0 ← ARG0  
[They] wrote [an email] to *cancel* it.

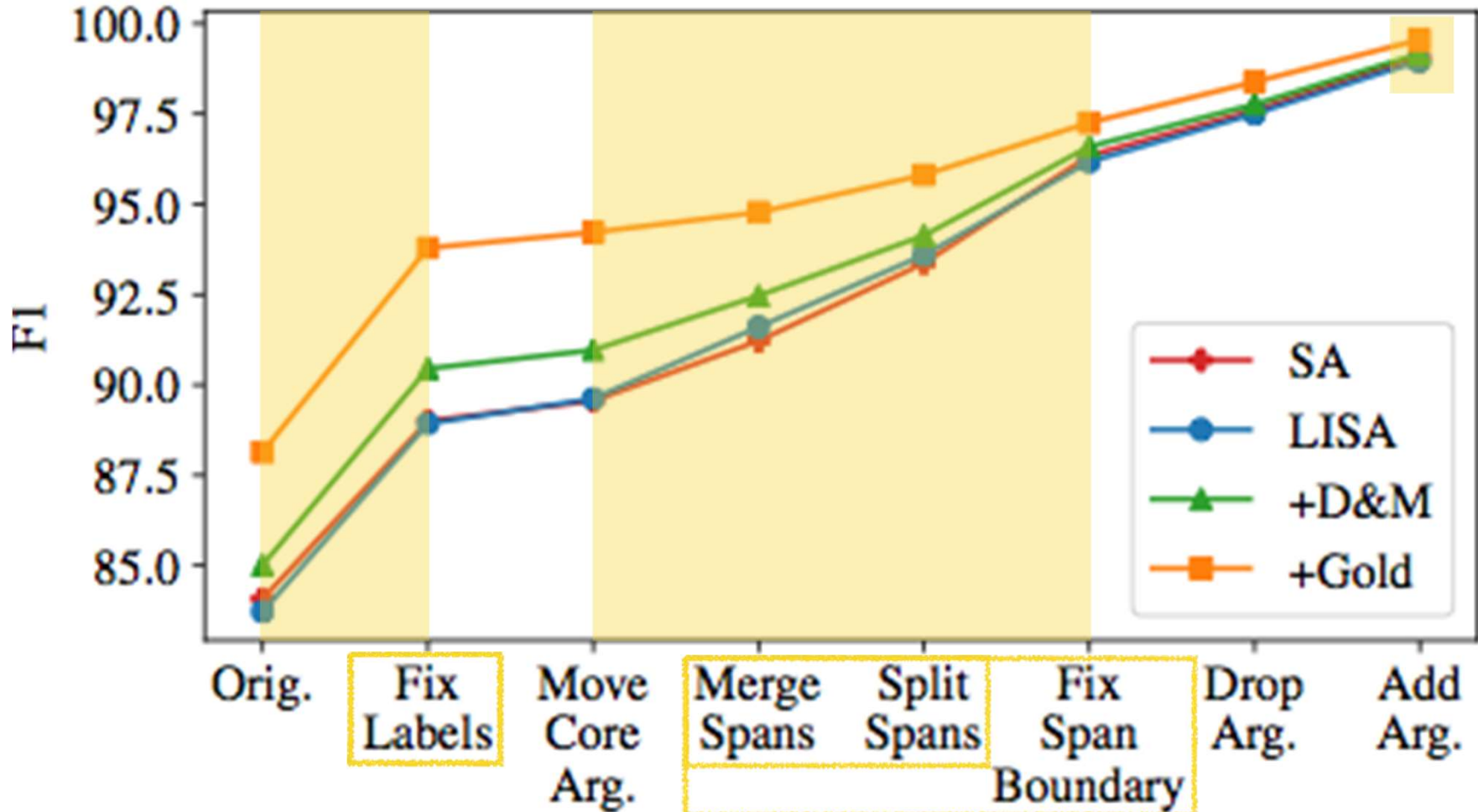
3) Split/  
Merge span:  
I *eat* [pasta with delight].  
ARG1  
ARG1 ARGM-MNR  
[pasta] [with delight]

ARG1 ARGM-MNR  
I *eat* [pasta] [with broccoli].  
ARG1  
[pasta with broccoli]

4) Fix span boundary:  
ARG1  
[“No broccoli”,] I said.  
[“No broccoli”],

5) Drop/  
add arg:  
Hesitantly, they *declined* to elaborate [on that matter].  
+ ARGM-MNR  
[Hesitantly], they *declined* to elaborate on that matter.  
~~ARG1~~

# Experimental results: Analysis



boundary mistakes

# Summary

- **LISA**: Multi-task learning + multi-head self attention trained to attend to syntactic parents
  - Achieves state-of-the-art F1 on PropBank SRL
  - Linguistic structure improves generalization
  - Fast: encodes sequence *only once* to predict predicates, parts-of-speech, labeled dependency parse, SRL

Slide reference: Emma Strubell  
on Linguistically-Informed Self-Attention for Semantic Role Labeling (best  
paper, EMNLP 2018)