# Learning a Part-of-Speech Tagger from Two Hours of Annotation

Dan Garrette and Jason Baldridge

University of Texas at Austin

# Low-Resource Languages

6,900 languages in the world

~30 have non-negligible quantities of data

No million-word corpus for any endangered language

[Maxwell and Hughes, 2006]
[Abney and Bird, 2010]

# Low-Resource Languages

Kinyarwanda

   Niger-Congo; morphologically-rich

Malagasy

   Austronesian; spoken in Madagascar

Also, English

# Low-Resource Languages

Supervised training is not an option.

We do semi-supervised training.

➡️ Annotate some data by hand

... cheaply

... like, in 2 hours

# Semi-Supervised Training

HMM with Expectation-Maximization (EM)

Need:

Large **raw** corpus ⬅ know how to get this

Tag dictionary ⬅ where is this from?

[Kupiec, 1992]
[Merialdo, 1994]

# Tag Dictionary

Most previous work:

Extract from a **large labeled corpus**

➡️ too **complete**

➡️ too **clean**

➡️ too **biased**

**unrealistic**

# A **Real** Tag Dictionary

# A **Real** Tag Dictionary

Extremely low coverage means most words are unknown
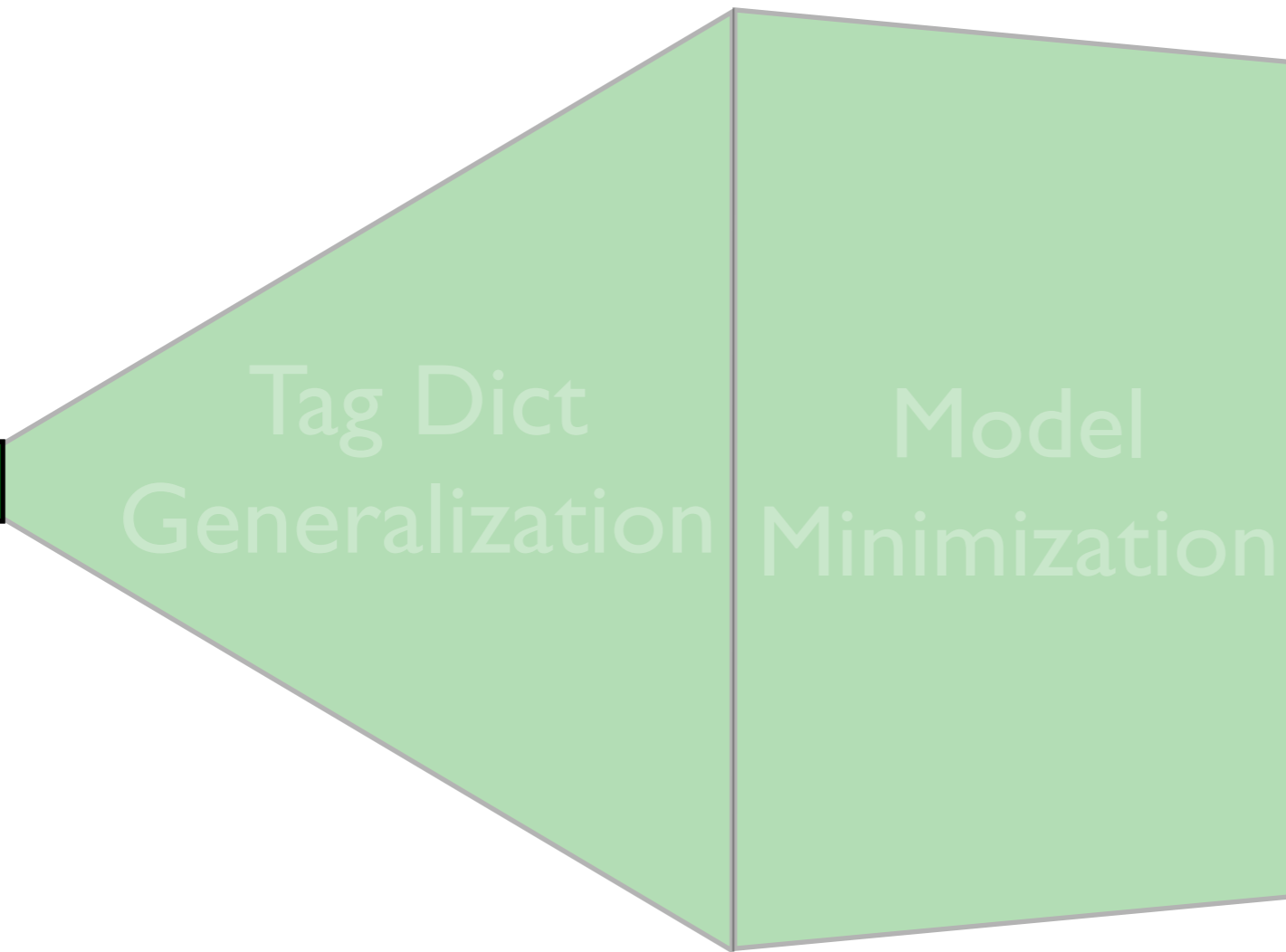
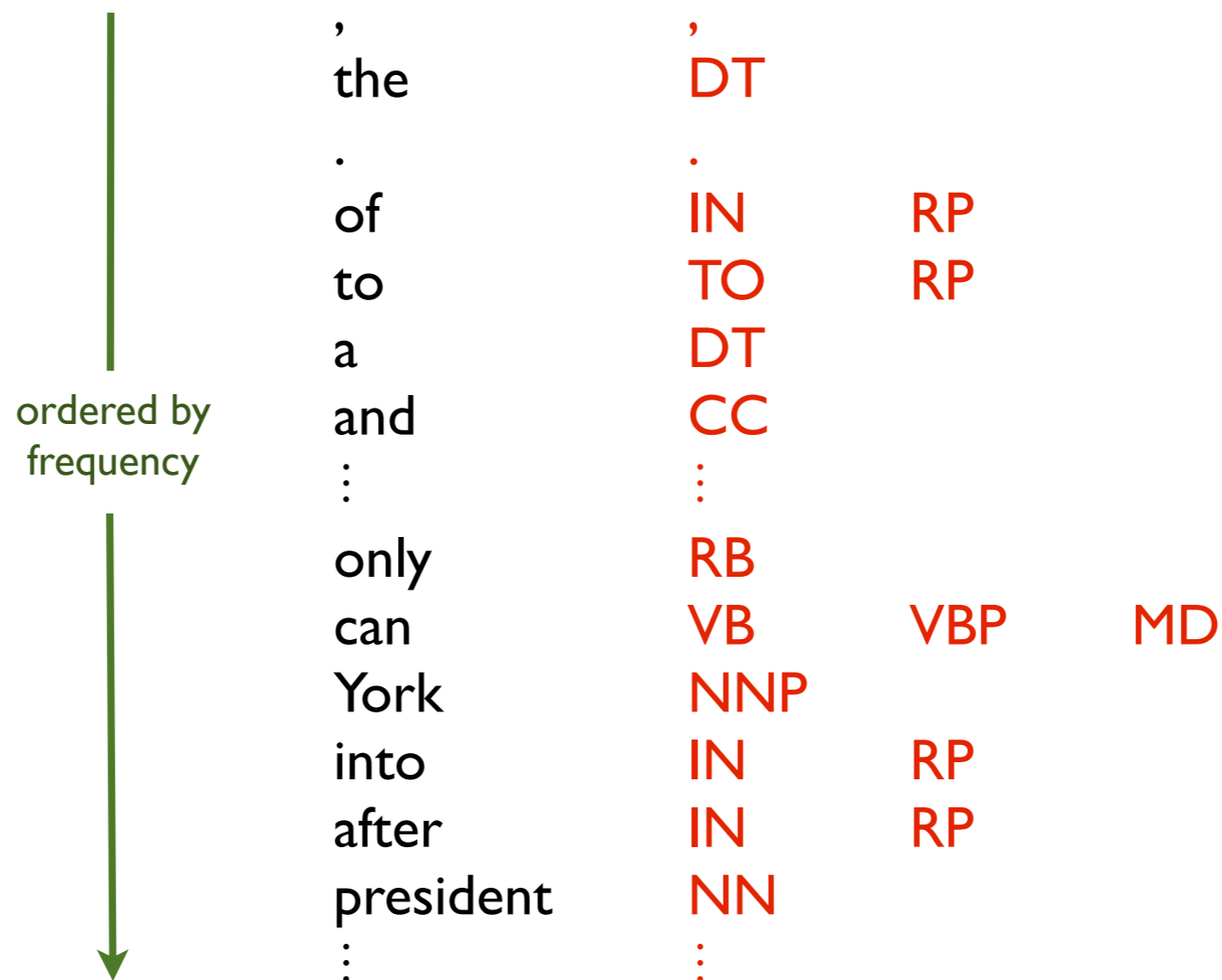$\Rightarrow$ **Bad for EM**  (poorly constrained)

# Our Approach



annotation → Tag Dict Generalization → Model Minimization → EM → HMM

cover the vocabulary    remove noise    train

# Collecting Annotations

Task #1 -- **2 hours** to create a **tag dictionary**

|  | word | tag | | |
|---|---|---|---|---|
| | , | , | | |
| | the | DT | | |
| | . | . | | |
| | of | IN | RP | |
| | to | TO | RP | |
| | a | DT | | |
| | and | CC | | |
| | ⋮ | ⋮ | | |
| | only | RB | | |
| | can | VB | VBP | MD |
| | York | NNP | | |
| | into | IN | RP | |
| | after | IN | RP | |
| | president | NN | | |
| | ⋮ | ⋮ | | |

ordered by frequency

# Collecting Annotations

Task #2 -- **2 hours** to annotate **full sentences**

| Pierre | Vinken | , | 61 | years | old | , | will | join | the | board | as | a | nonexecutive | director | Nov. | 29 | . |
|--------|--------|---|----|-------|-----|---|------|------|-----|-------|----|----|--------------|----------|------|----|---|
| NNP | NNP | , | CD | NNS | JJ | , | MD | VB | DT | NN | IN | DT | JJ | NN | NNP | CD | . |

| Mr. | Vinken | is | chairman | of | Elsevier | N.V. | , | the | Dutch | publishing | group | . |
|-----|--------|----|----------|----|----------|------|---|-----|-------|------------|-------|---|
| NNP | NNP | VB | NN | IN | NNP | NNP | , | DT | JJ | JJ | NN | . |

⋮

# Collecting Annotations

In 2 hours:

|                | # sent | # tok | # TD entries |
|----------------|--------|-------|--------------|
| Full Sentences | 90     | 1537  | **750**      |
| Tag Dict       |        |       | **1798**     |

(for Kinyarwanda)

# Our Approach

# Tag Dict Generalization

These annotations are too sparse!

Generalize to the entire vocabulary

# Tag Dict Generalization

Haghighi and Klein (2006) do this with
a vector space.

<span style="color:red">We don't have enough raw data</span>

Das and Petrov (2011) do this with
a parallel corpus.

<span style="color:red">We don't have a parallel corpus</span>

# Tag Dict Generalization

Our strategy: Label Propagation

- **Connect** annotations to raw corpus tokens

- Push tag labels to **entire corpus**

[Talukdar and Crammer. 2009]

# Tag Dict Generalization

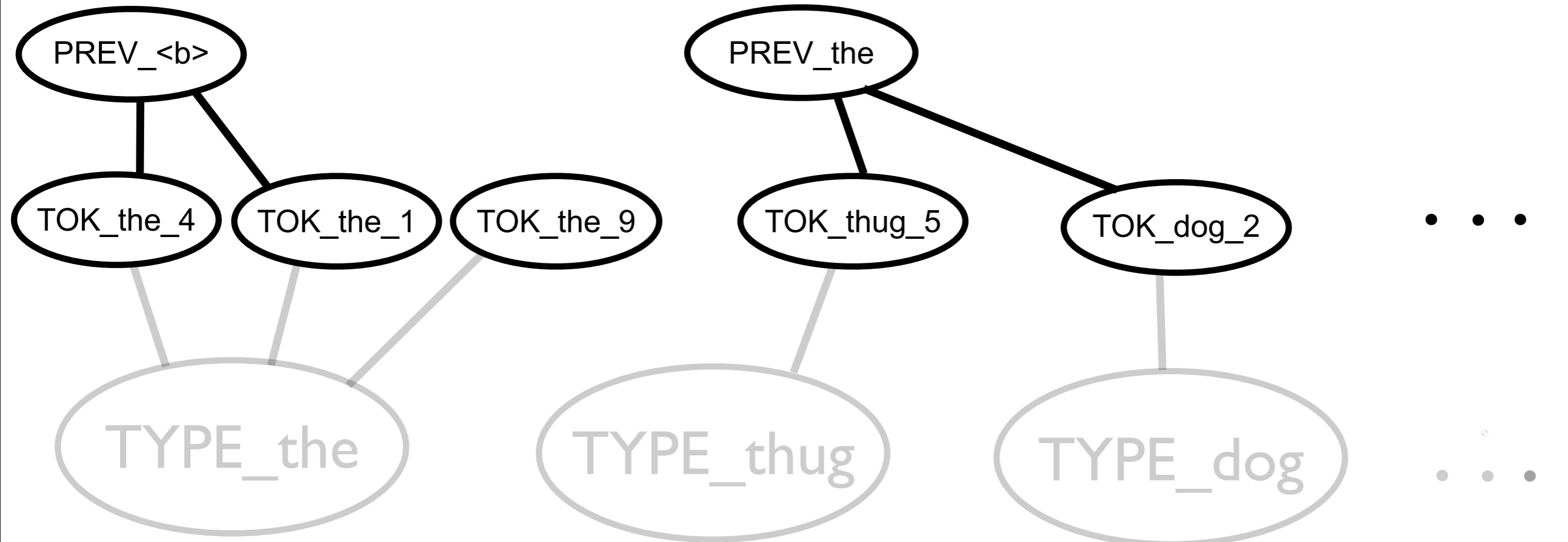Annotations

Raw Corpus

TYPE_the    TYPE_thug    TYPE_dog    . . .

# Tag Dict Generalization

Annotations

the$_4$  thug$_5$  walks$_6$

Raw Corpus

( TOK_the_4 )  ( TOK_thug_5 )  ( TOK_walks_6 )

( TYPE_the )  ( TYPE_thug )  ( TYPE_dog )  . . .

# Tag Dict Generalization

Annotations

Raw Corpus

the$_4$    thug$_5$    walks$_6$

TOK_the_4    TOK_the_1    TOK_the_9    TOK_thug_5    TOK_dog_2    • • •

TYPE_the    TYPE_thug    TYPE_dog    • • •

# Tag Dict Generalization

# Tag Dict Generalization

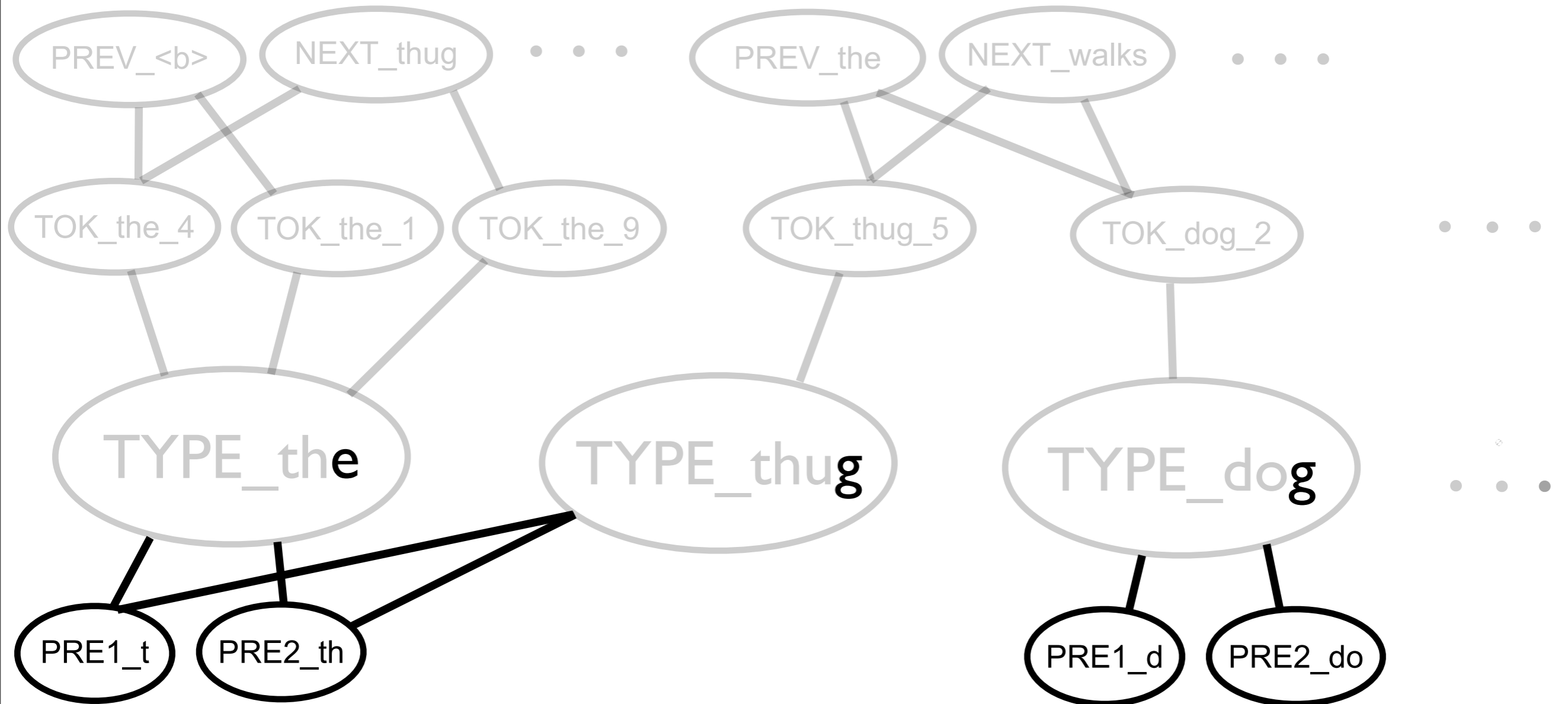Annotations

the$_4$     thug$_5$     walks$_6$

Raw Corpus

# Tag Dict Generalization

Annotations

the$_4$    thug$_5$    walks$_6$

Raw Corpus

NEXT_thug    NEXT_walks    NEXT_<b>

PREV_<b>    PREV_the

TOK_the_4    TOK_the_1    TOK_the_9    TOK_thug_5    TOK_dog_2    • • •

TYPE_the    TYPE_thug    TYPE_dog    • • •

# Tag Dict Generalization

Annotations

the$_4$    thug$_5$    walks$_6$

Raw Corpus

PREV_<b>  NEXT_thug  •  •  •  PREV_the  NEXT_walks  •  •  •

TOK_the_4  TOK_the_1  TOK_the_9  TOK_thug_5  TOK_dog_2  •  •  •

TYPE_the  TYPE_thug  TYPE_dog  •  •  •

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization
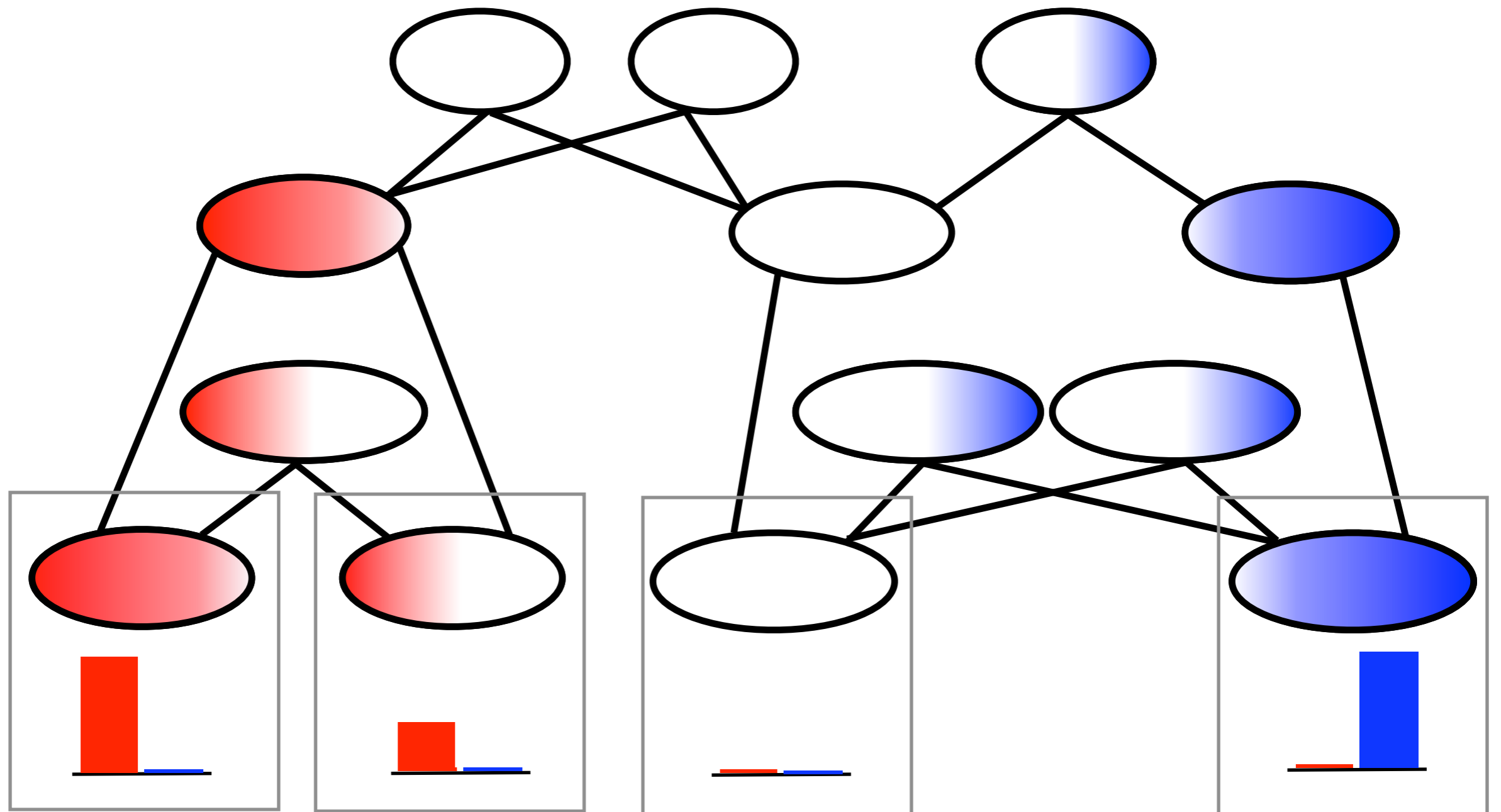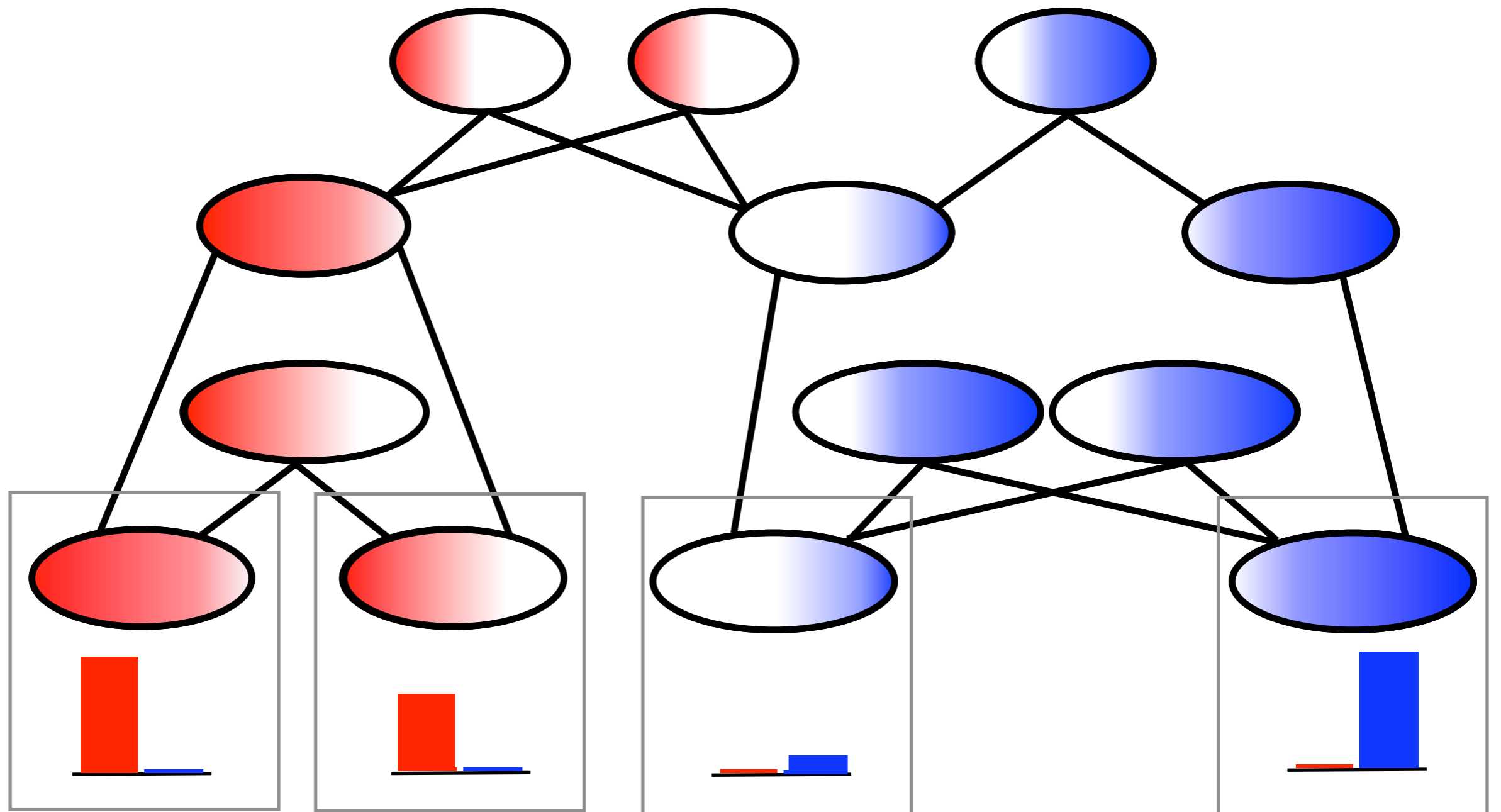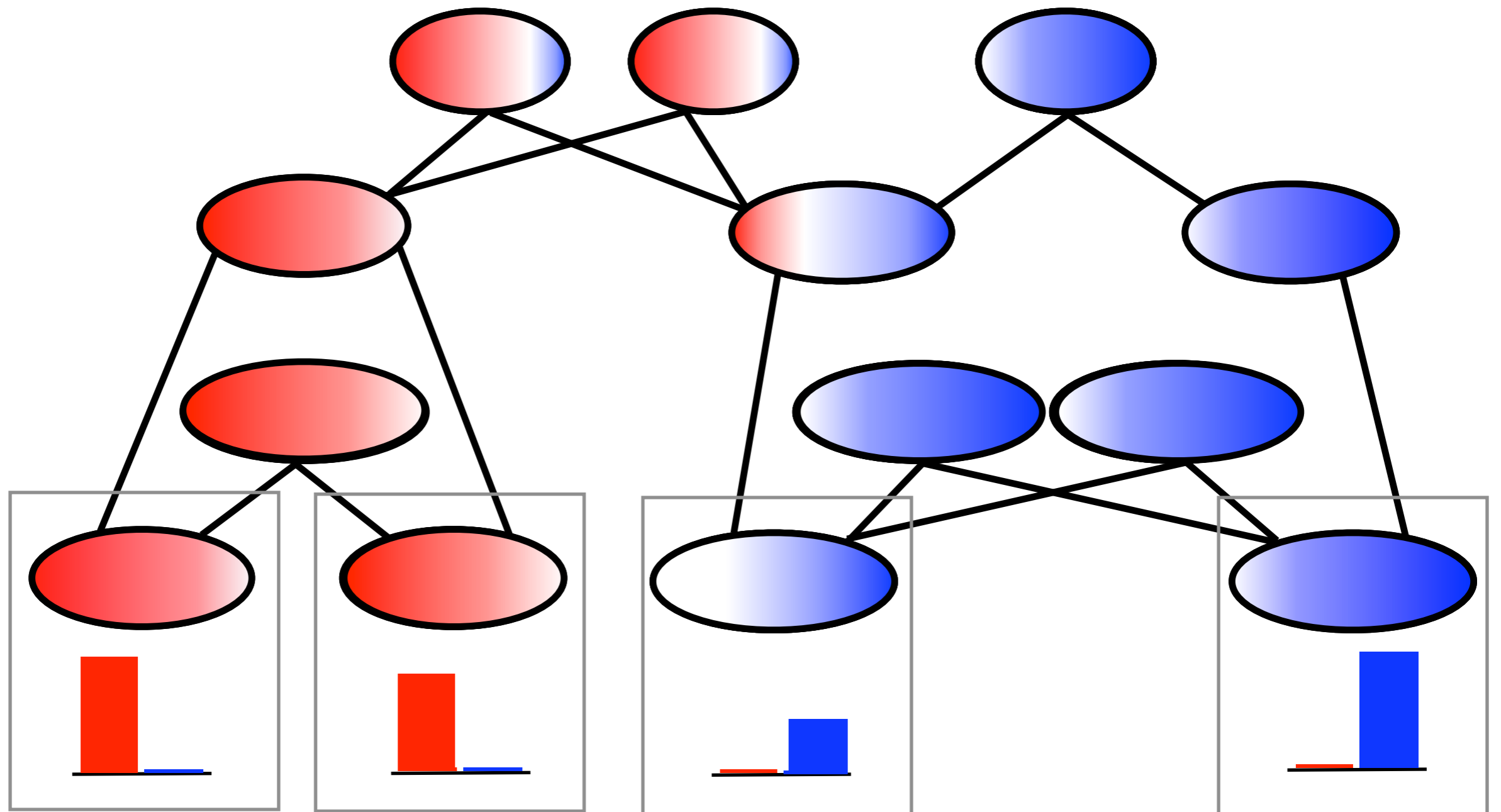
# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization
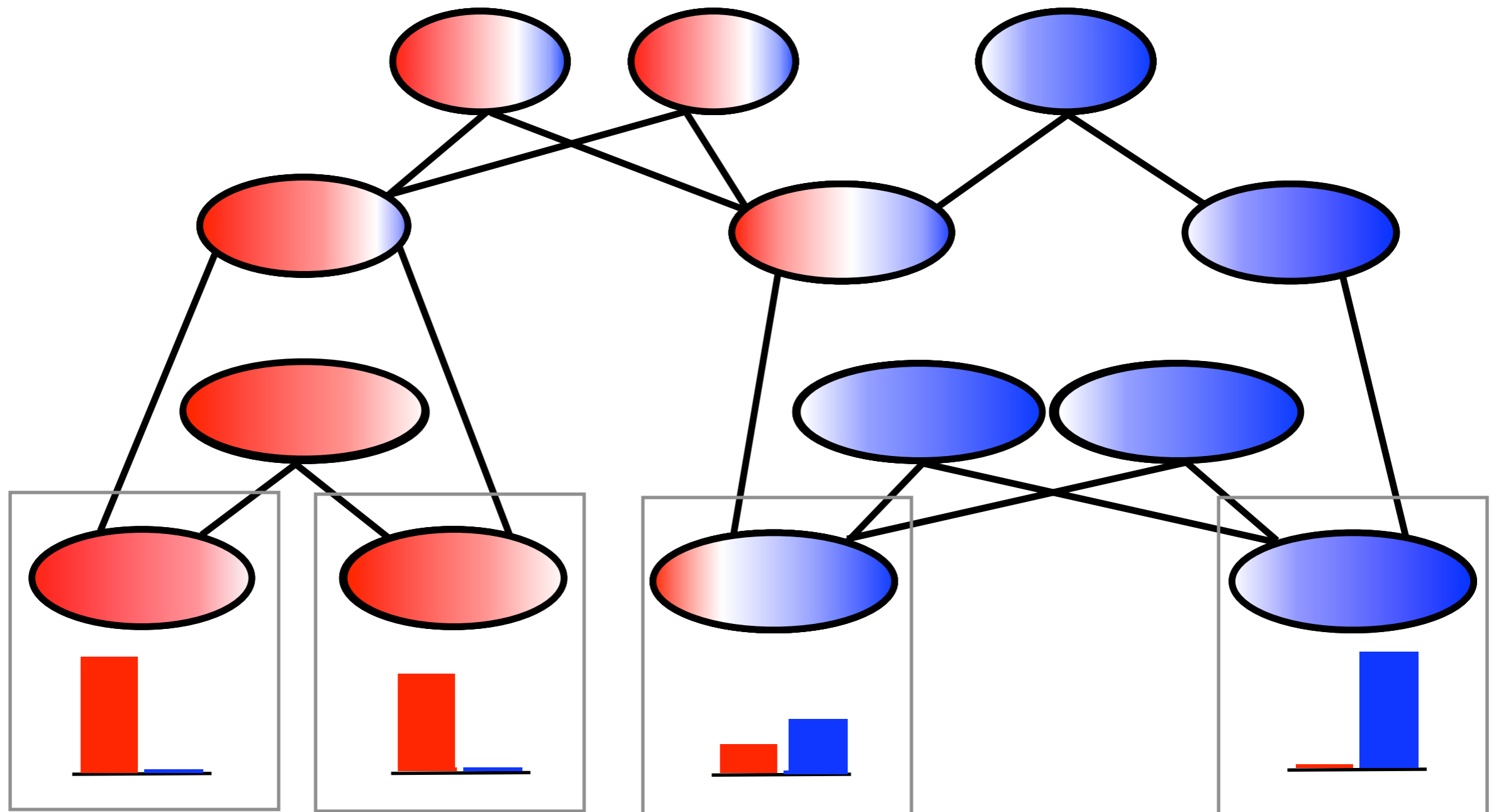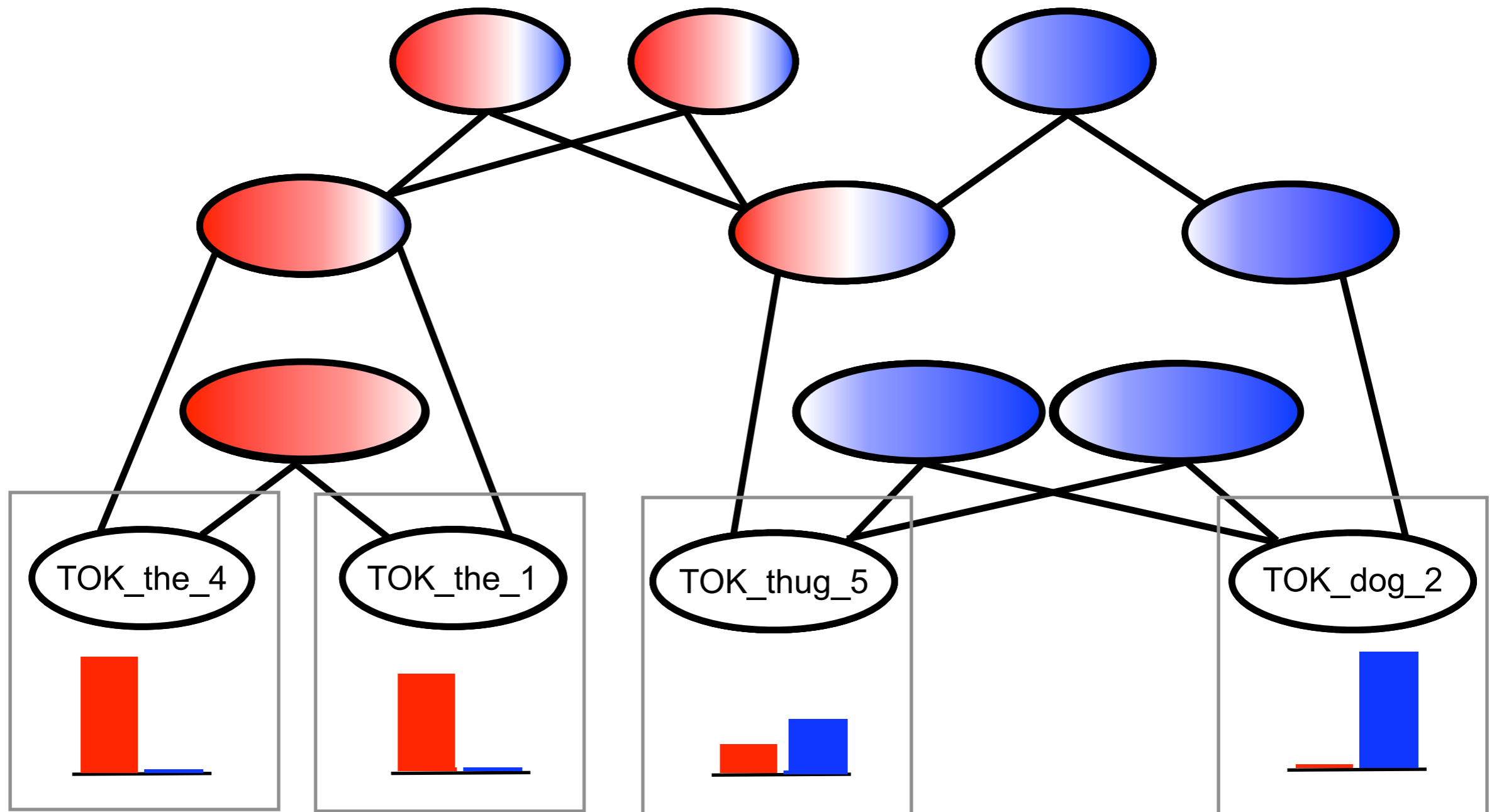
# Tag Dict Generalization

Type Annotations
the
dog

PRE2_th  PRE1_t  SUF1_g

TYPE_**DT**_the  TYPE_thug  TYPE_**NN**_dog

PREV_<b>  PREV_the  NEXT_walks

TOK_the_4  TOK_the_1  TOK_thug_5  TOK_dog_2

Token Annotations
the dog walks
**DT** **NN** **VBZ**

# Tag Dict Generalization



Type Annotations
the
dog

Token Annotations
the dog walks

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization
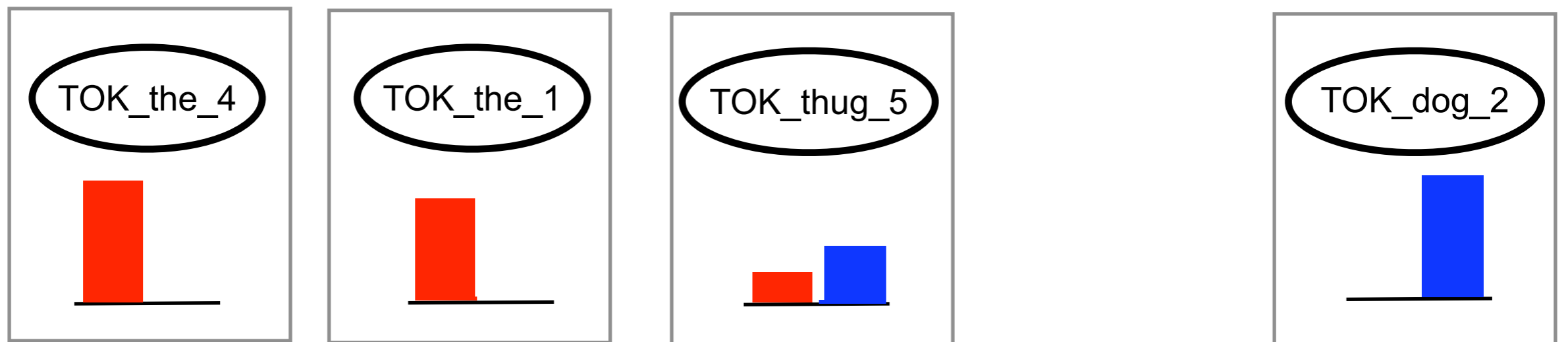
# Tag Dict Generalization
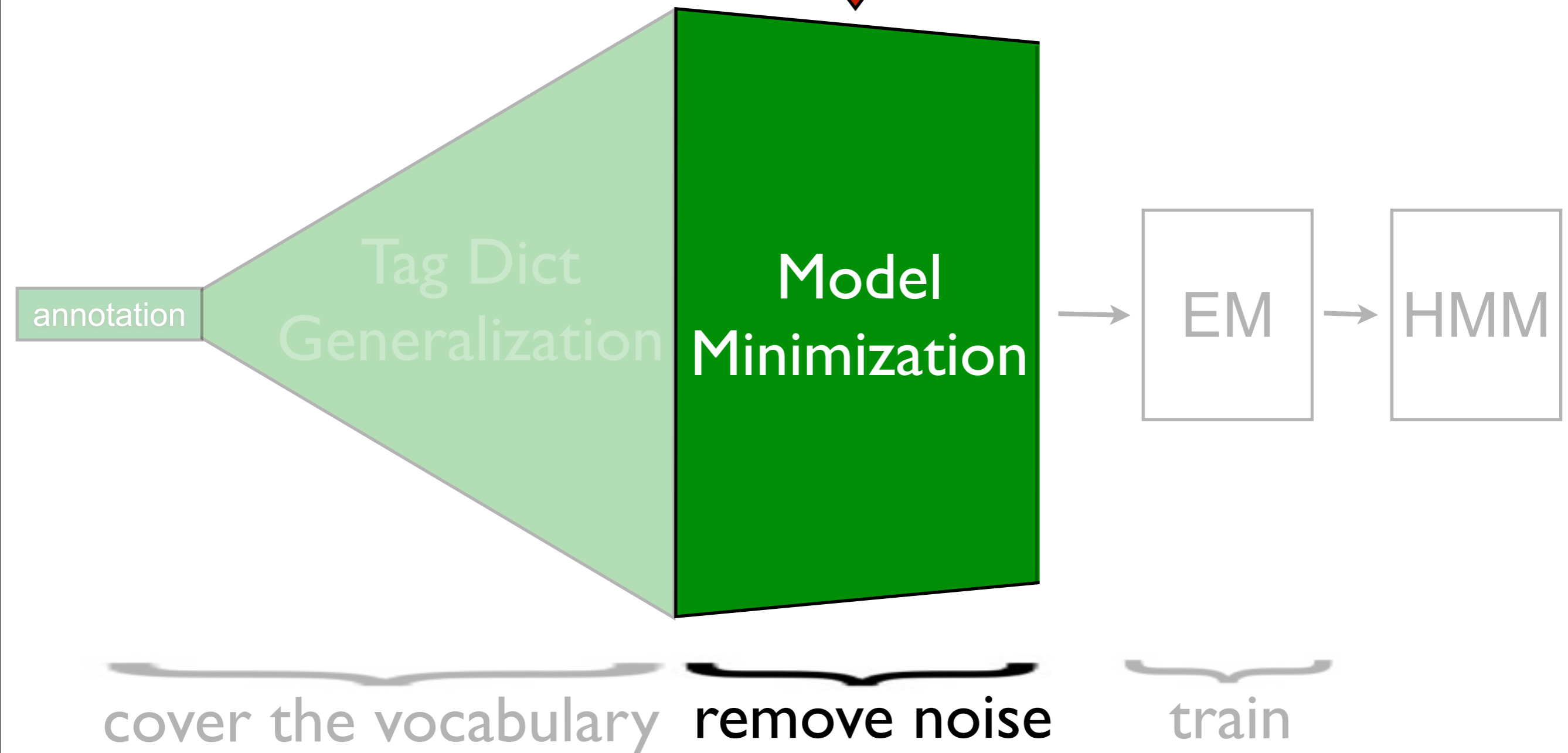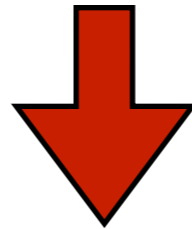
# Tag Dict Generalization

# Tag Dict Generalization

Result:

- a tag distribution on every token (soft tagging)
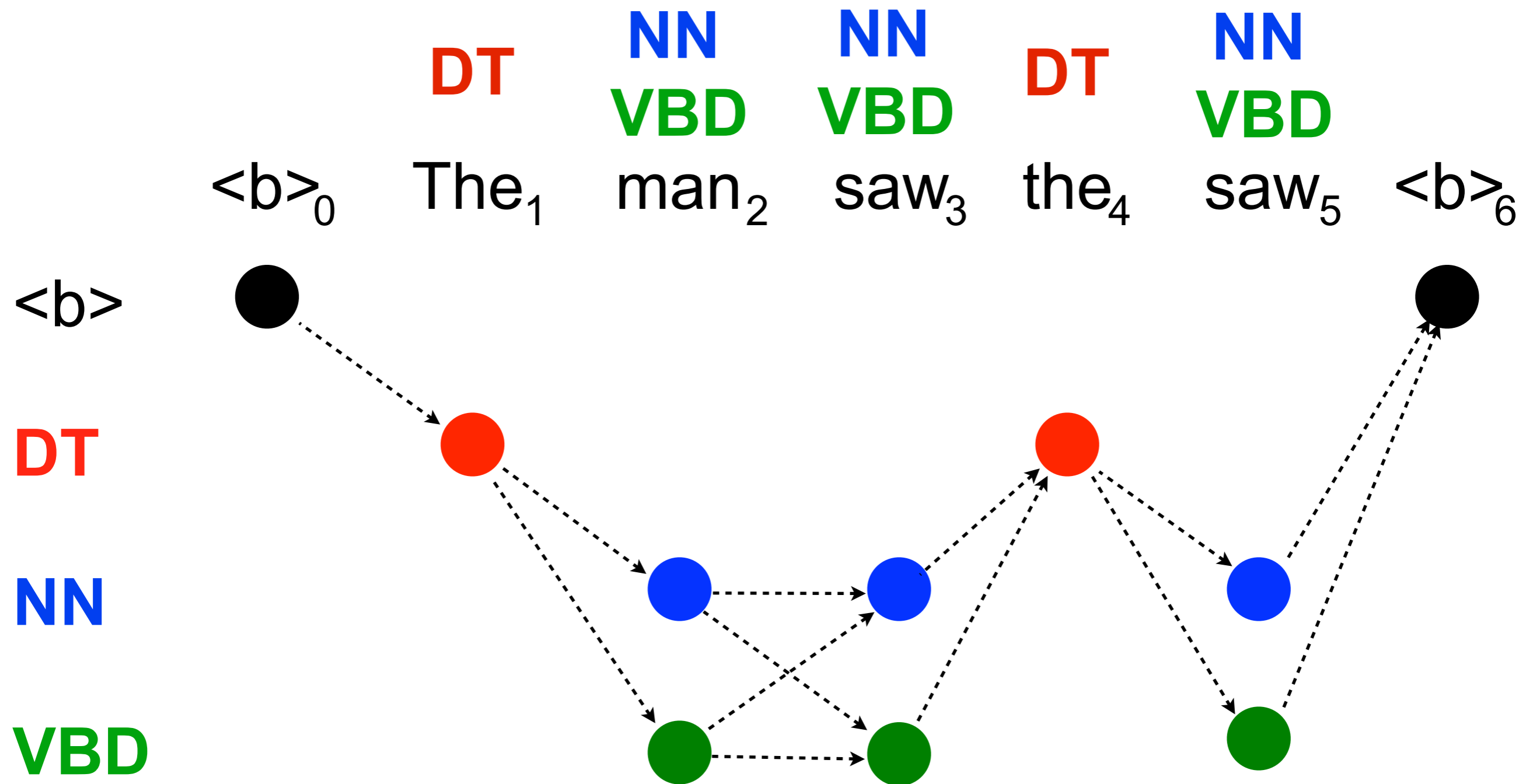
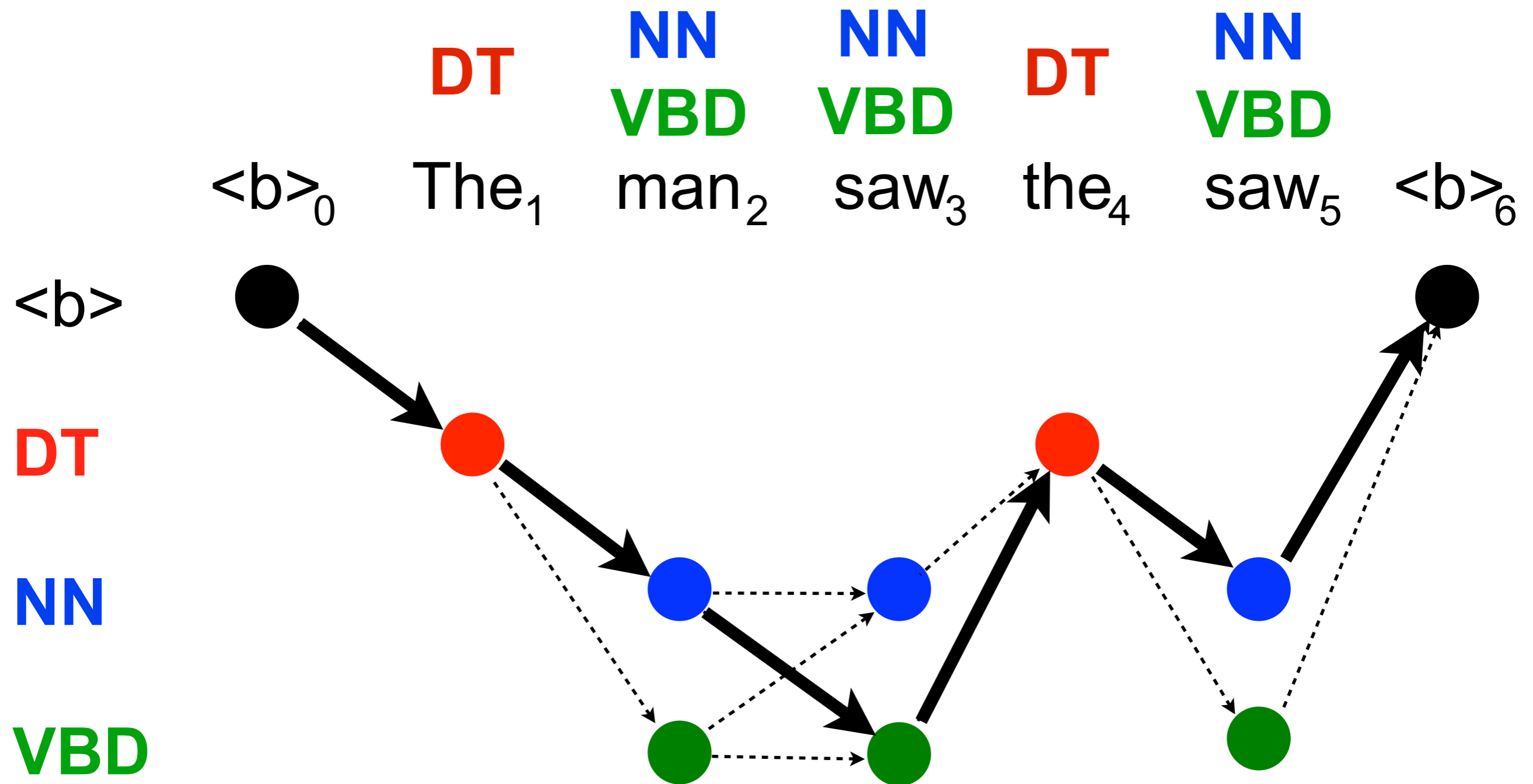- an expanded tag dictionary (non-zero tags)

# Our Approach

annotation

Tag Dict Generalization

**Model Minimization**

EM → HMM

cover the vocabulary | remove noise | train

# Model Minimization

- Induce a cleaner hard tagging from a noisy soft tagging.

- Greedily seek the minimal set of tag bigrams that describe the raw corpus.

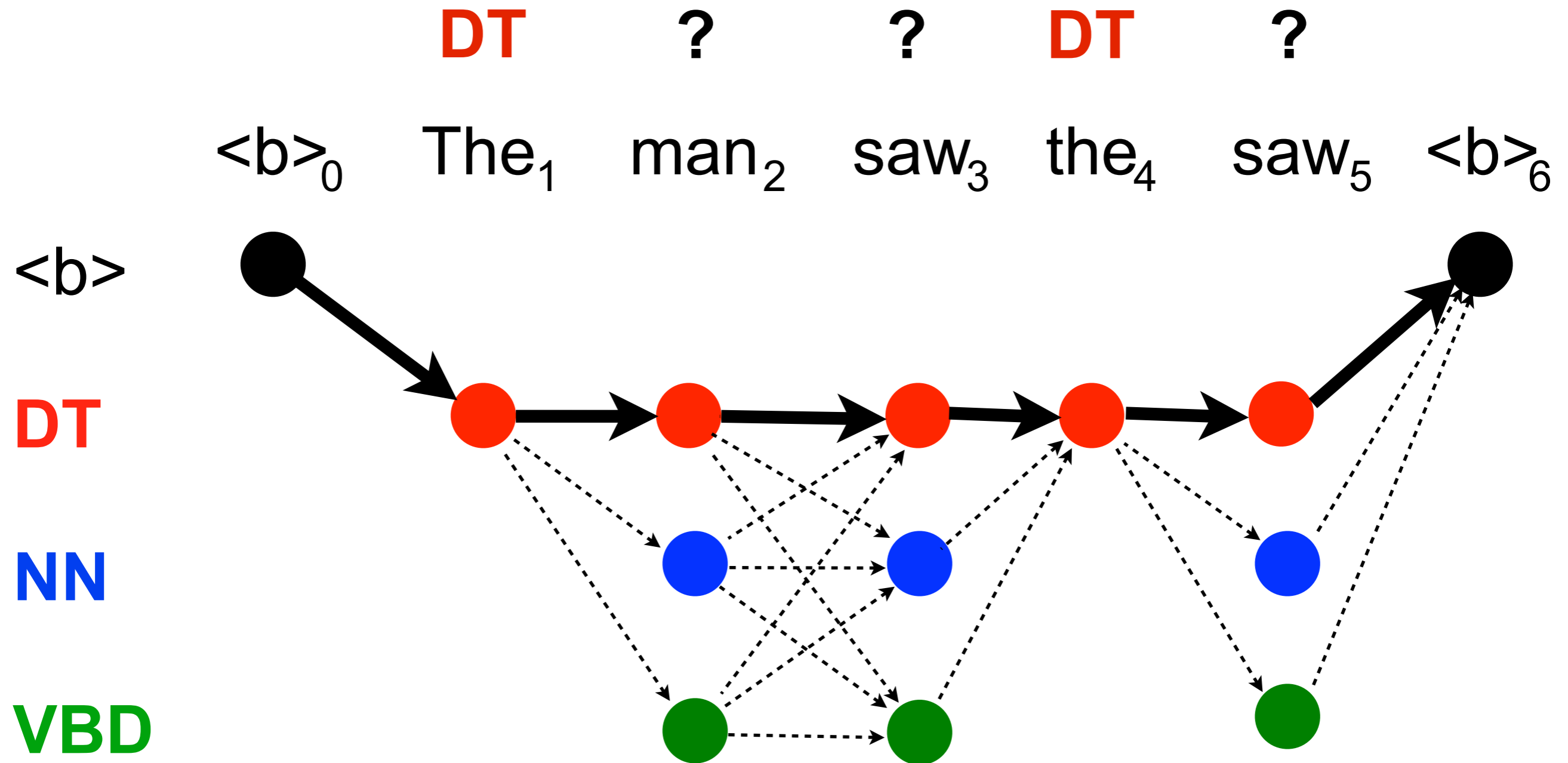[Ravi et al., 2010; Garrette and Baldridge, 2012]

# Model Minimization
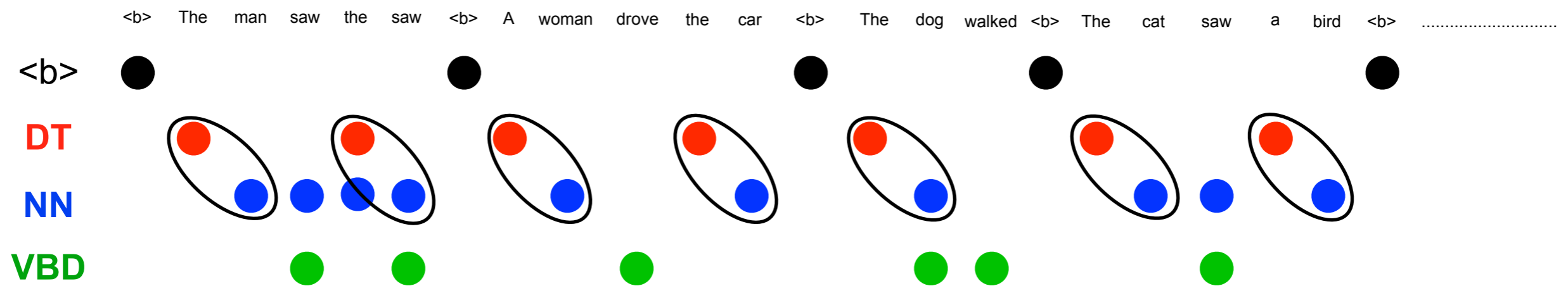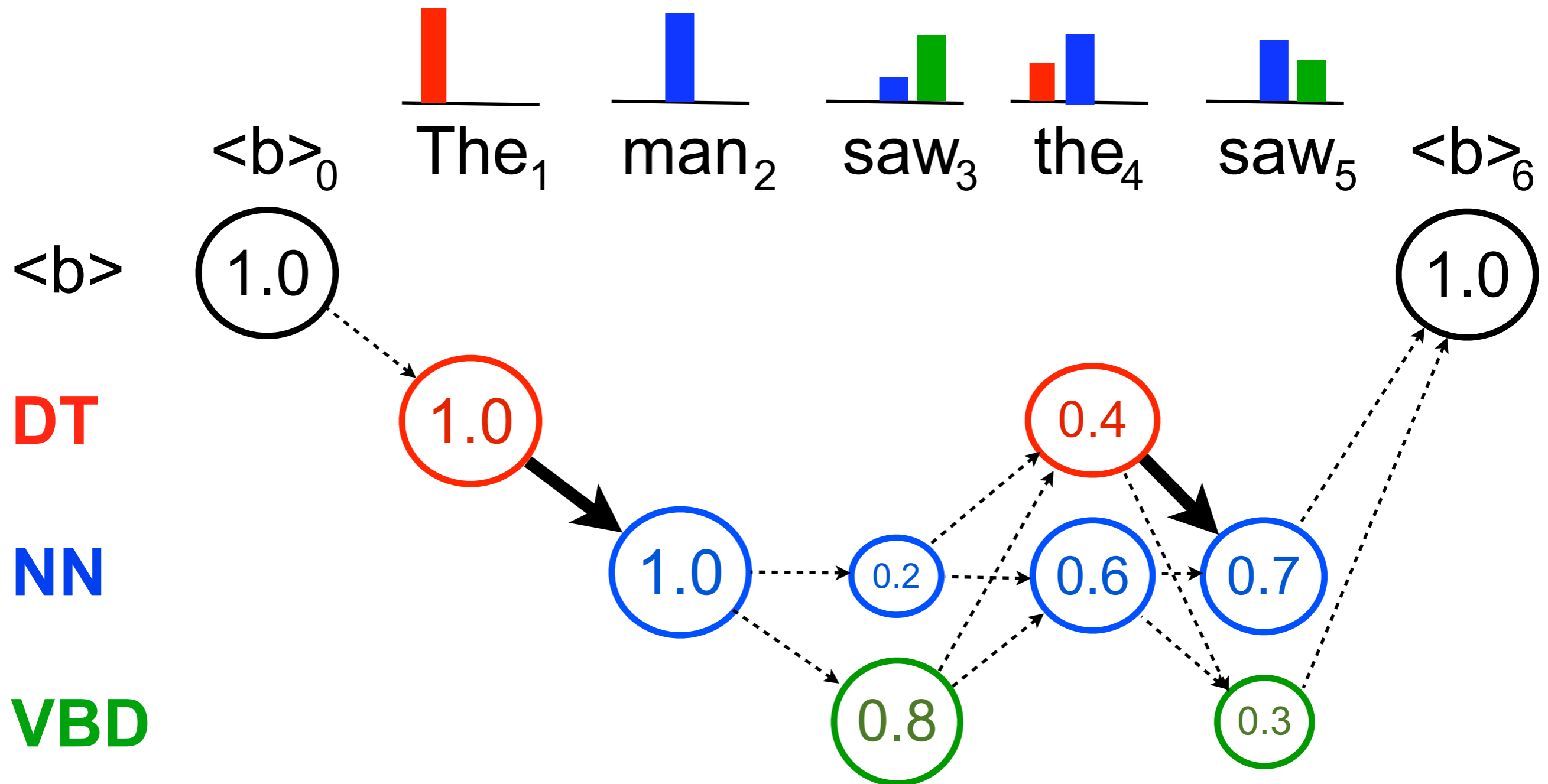
# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

The$_1$ man$_2$ saw$_3$ the$_4$ saw$_5$ &lt;b&gt;$_6$

&lt;b&gt;$_0$

&lt;b&gt; 1.0

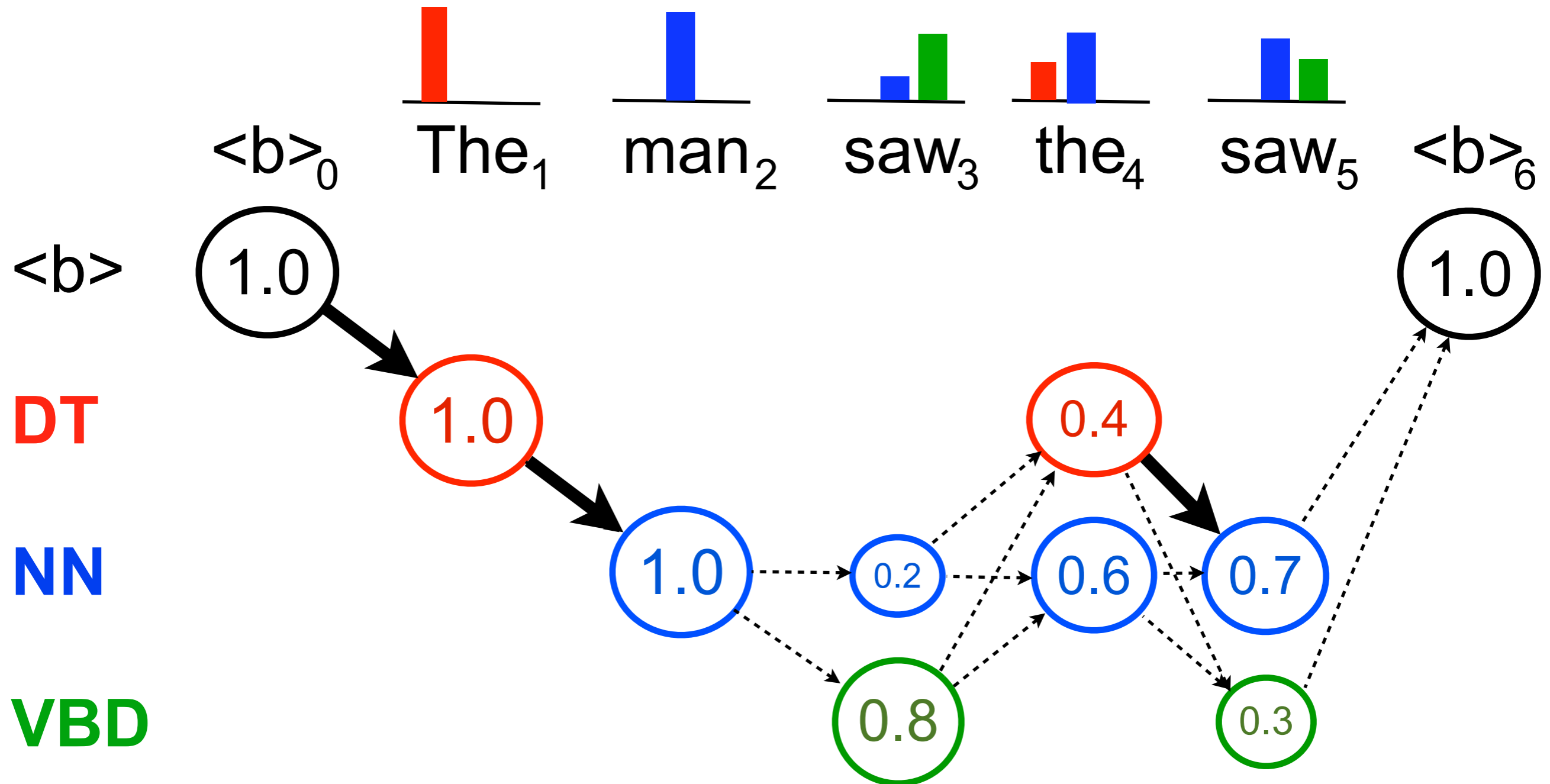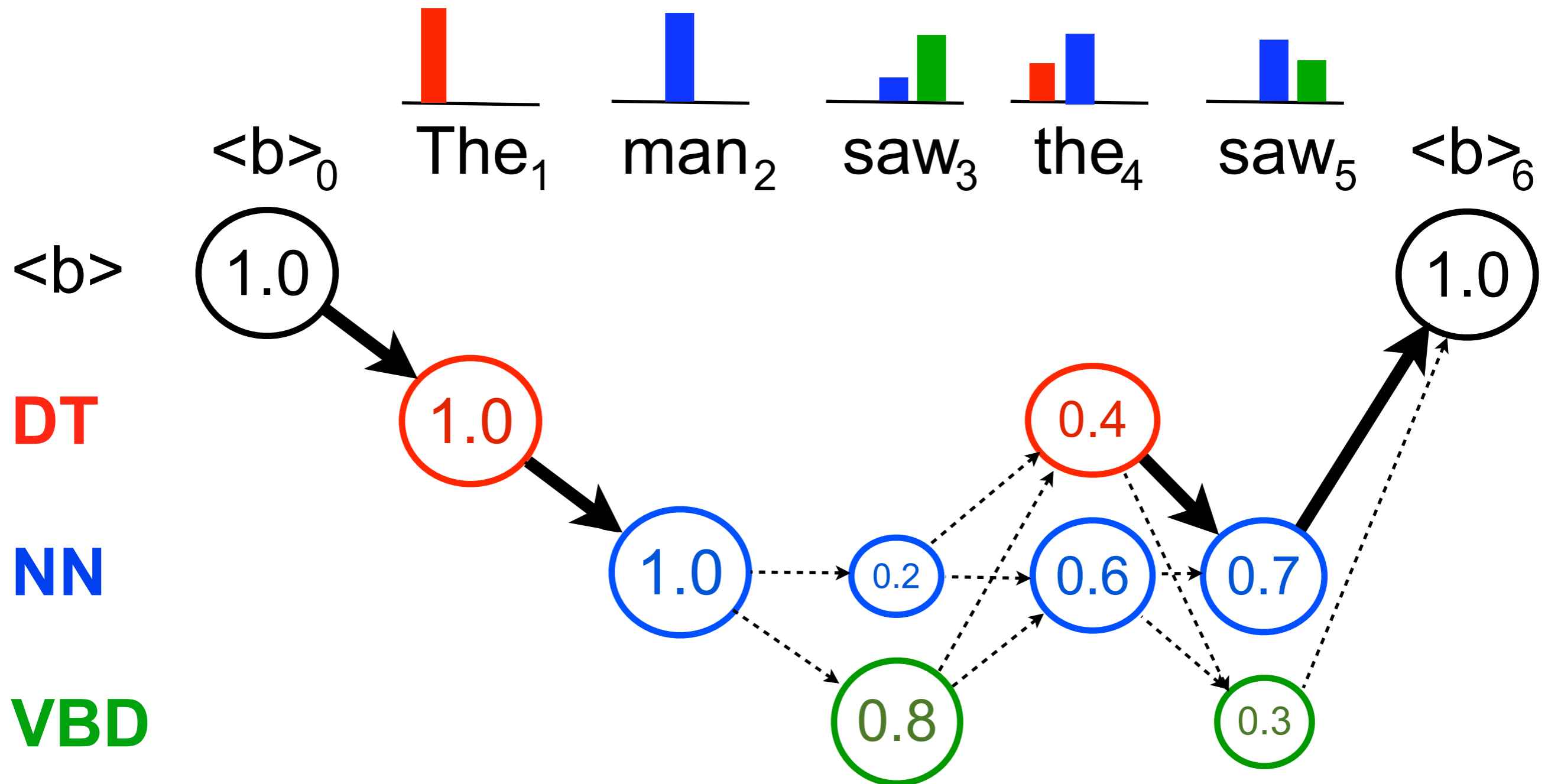DT 1.0 0.4

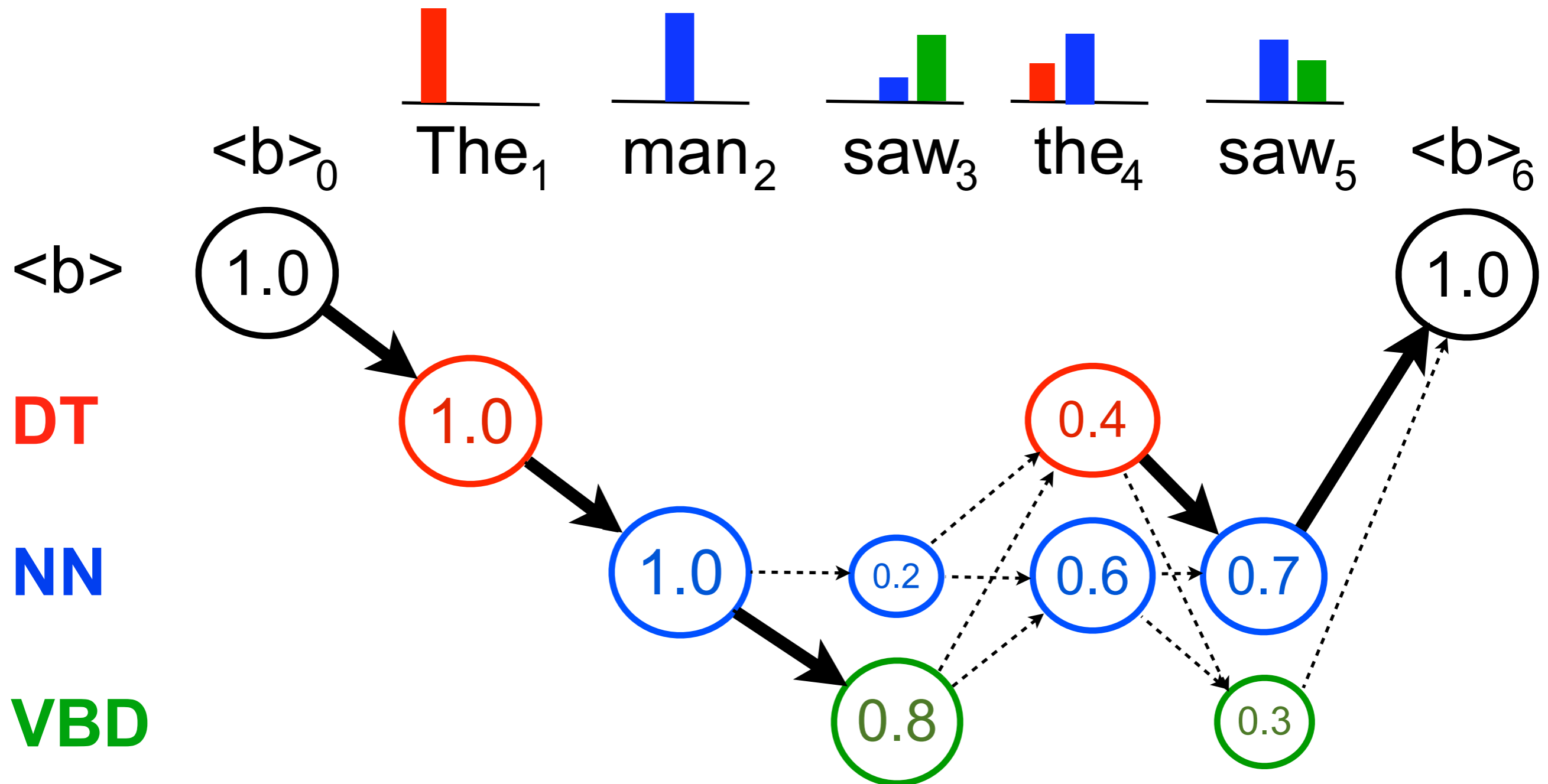NN 1.0 0.2 0.6 0.7

VBD 0.8 0.3

1.0
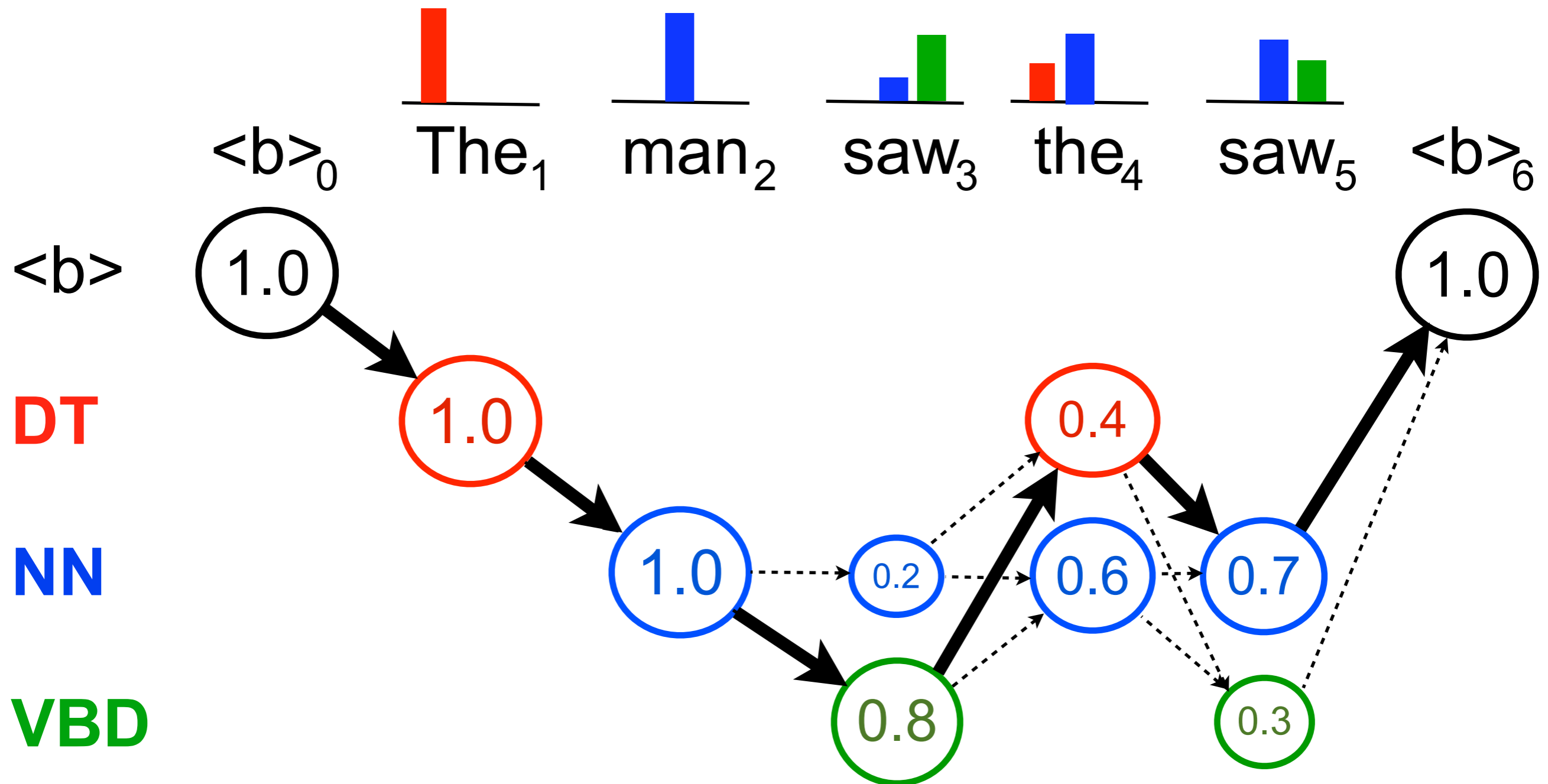
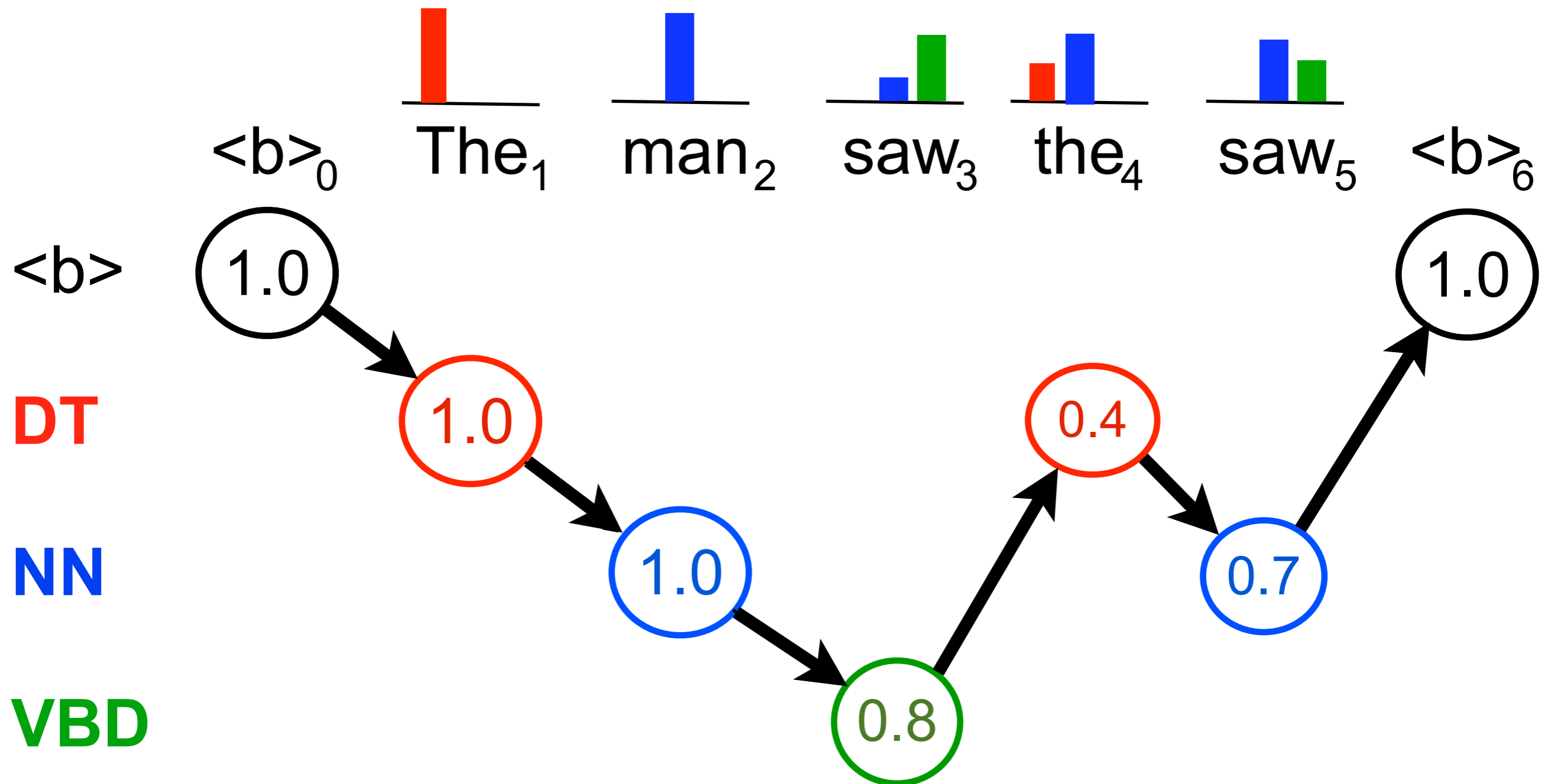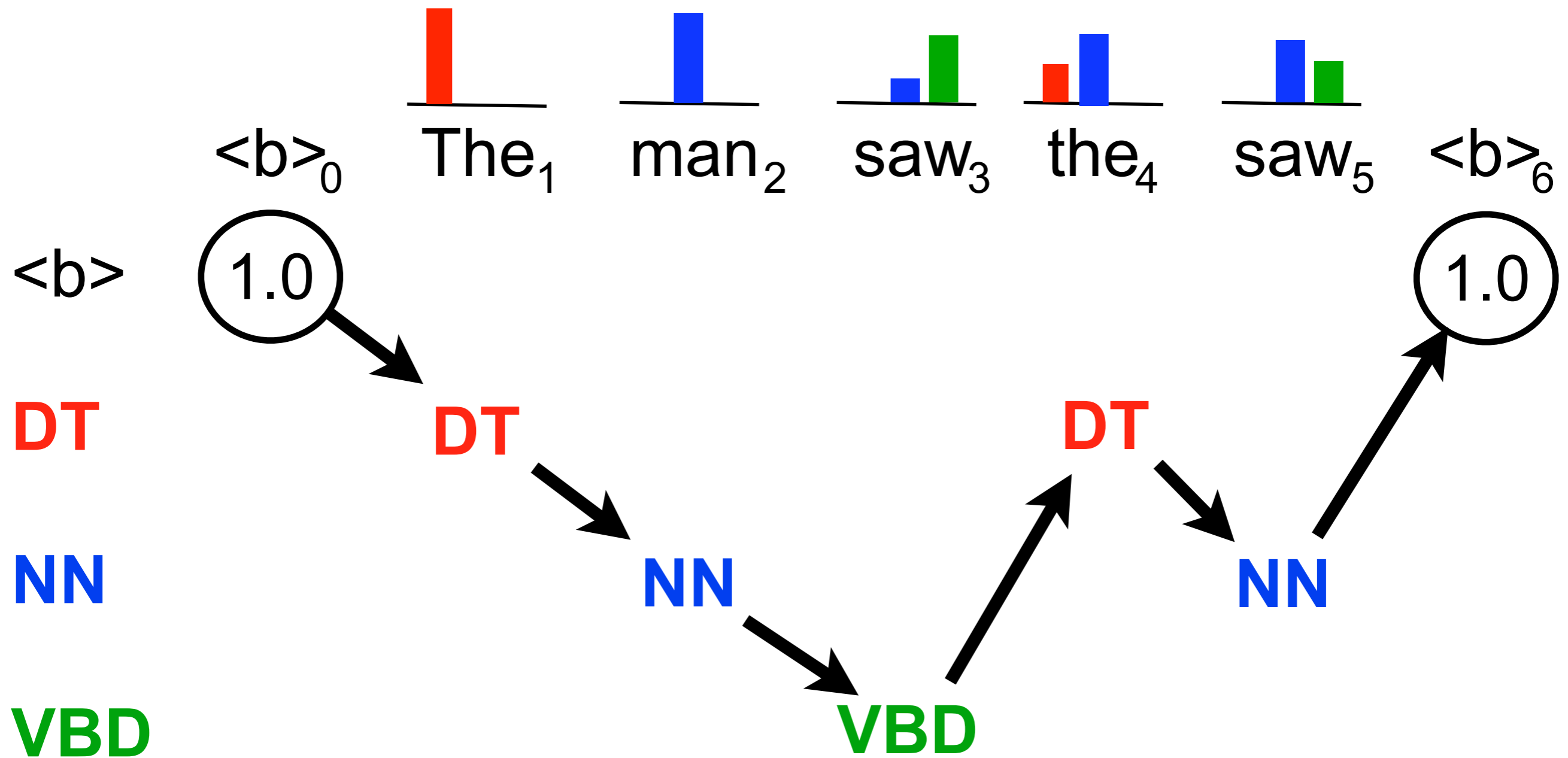# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

$\langle b\rangle_0$ The$_1$ man$_2$ saw$_3$ the$_4$ saw$_5$ $\langle b\rangle_6$

**DT** **NN** **VBD** **DT** **NN**

# Our Approach

# Model Minimization

# Unknown Accuracy

| | English | Kinyarwanda | Malagasy |
|---|---|---|---|
| Tokens — EM only | 43 | 32 | 39 |
| Types — EM only | 38 | 32 | 48 |
| Tokens — + Our approach | 61 | 58 | 53 |
| Types — + Our approach | 61 | 70 | 60 |

**Tokens**
- EM only
- + Our approach

**Types**
- EM only
- + Our approach

# Unknown Accuracy

Remember: Very high unknown rates.

Especially for morphological-rich Kinyarwanda.

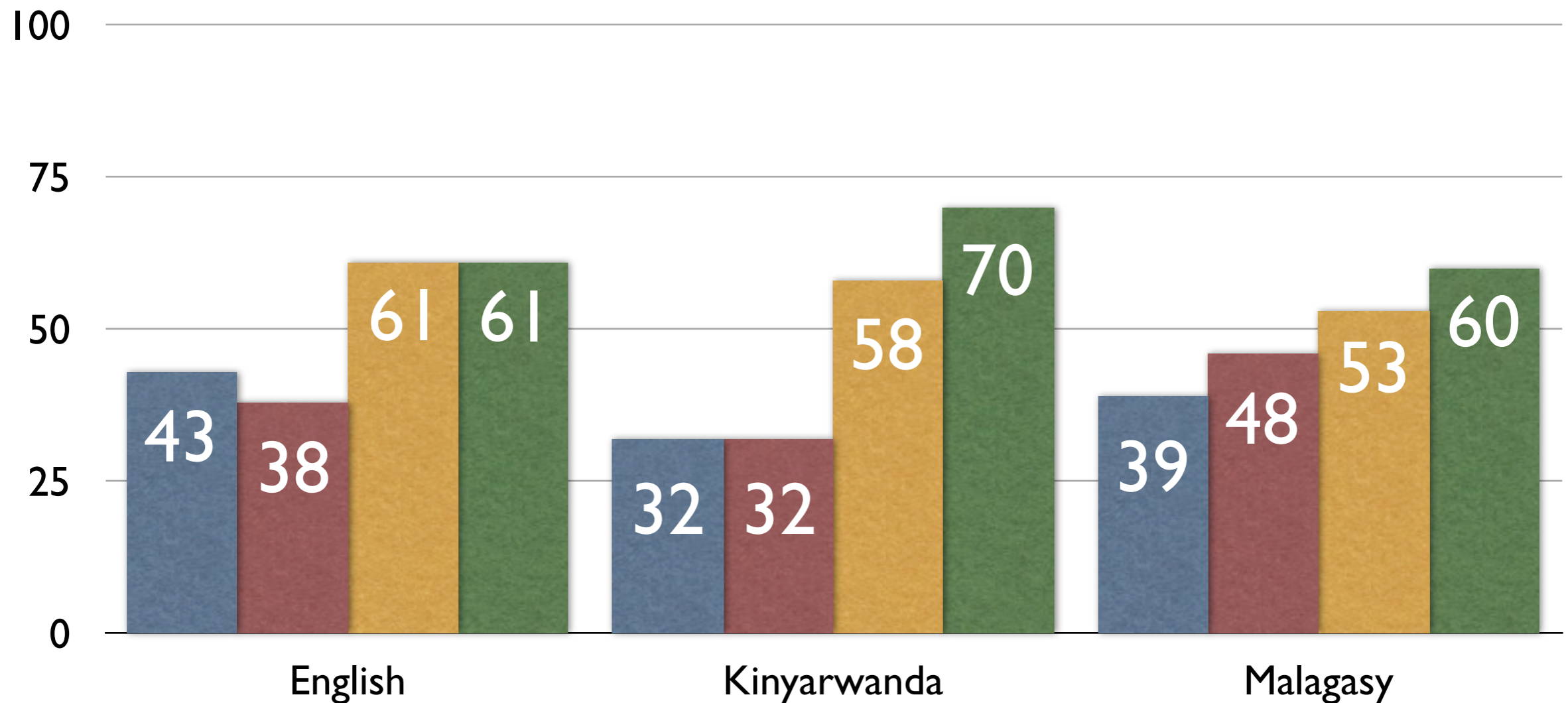# Conclusion

- Developed a semi-supervised approach to learn a tagger from realistically minimal input.

- Currently being used for further low-resource research (e.g. unsupervised dependency parsing).

# ACL Preview

- Learning curves for annotation time

- Mixed types and tokens under fixed time constraints

- Morphological transducers

- **90%** accuracy on full 45 tag English Penn Treebank with **4 hours** of data

# Software Available

**Train your own** low-resource taggers.

Or use our Kinyarwanda and Malagasy models.

Open source: link on my website or in the paper.