# Learning a Part-of-Speech Tagger from Minimal Annotation

Dan Garrette

University of Texas at Austin

# Low-Resource Languages

Supervised training is not an option.

# Low-Resource Languages

Supervised training is not an option.

We do semi-supervised training.

# Low-Resource Languages

Supervised training is not an option.

We do semi-supervised training.

→ Annotate some data by hand

# Low-Resource Languages

Supervised training is not an option.

We do semi-supervised training.

→ Annotate some data by hand

... cheaply

# Semi-Supervised Training

[Kupiec, 1992]
[Merialdo, 1994]

# Semi-Supervised Training

HMM with Expectation-Maximization (EM)

[Kupiec, 1992]
[Merialdo, 1994]

# Semi-Supervised Training

HMM with Expectation-Maximization (EM)

Need:

[Kupiec, 1992]
[Merialdo, 1994]

# Semi-Supervised Training
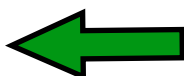
HMM with Expectation-Maximization (EM)

Need:

Large **raw** corpus

[Kupiec, 1992]
[Merialdo, 1994]

# Semi-Supervised Training

HMM with Expectation-Maximization (EM)

Need:

Large **raw** corpus

Tag dictionary

[Kupiec, 1992]
[Merialdo, 1994]

# Semi-Supervised Training

HMM with Expectation-Maximization (EM)
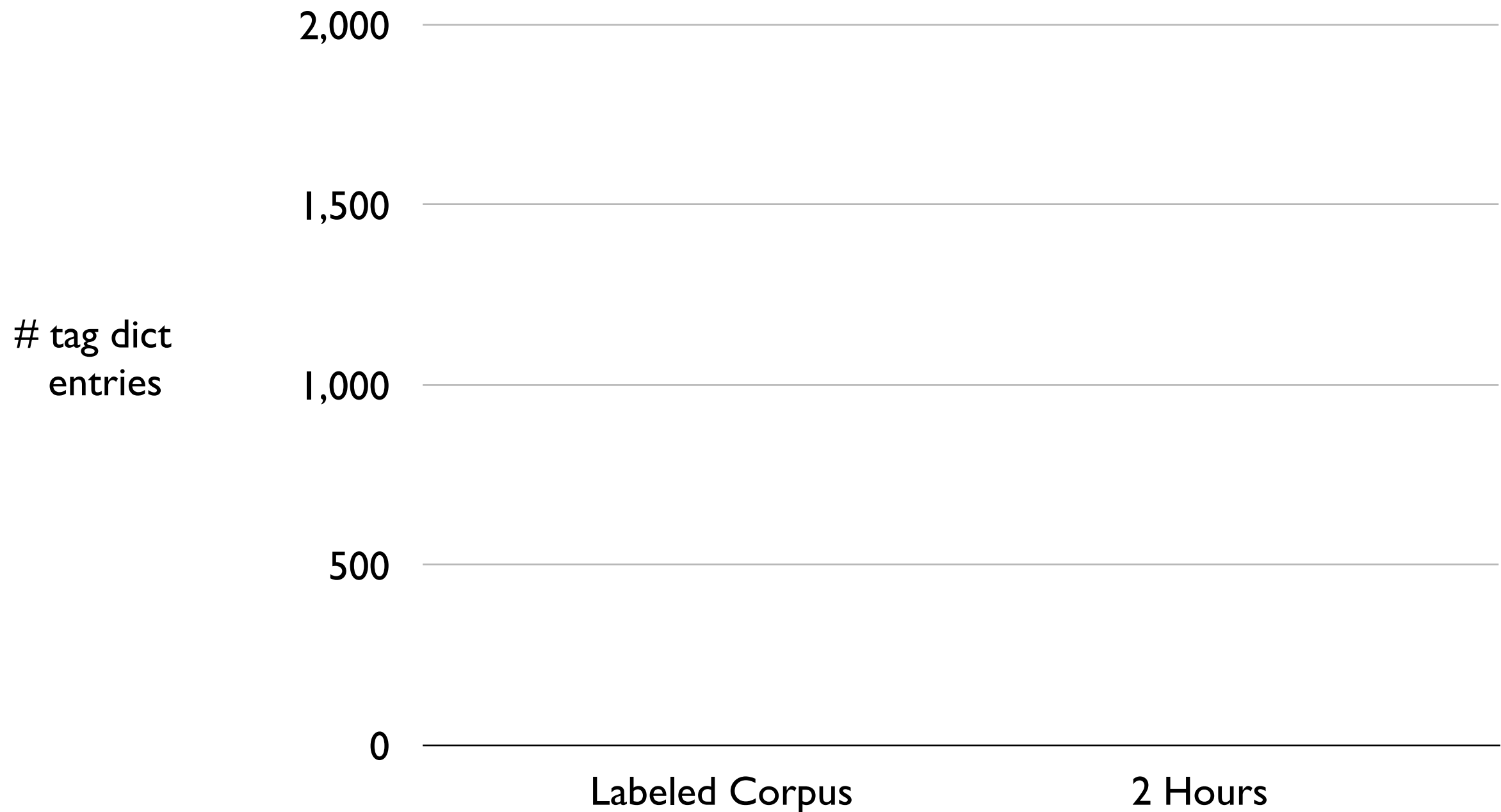
Need:

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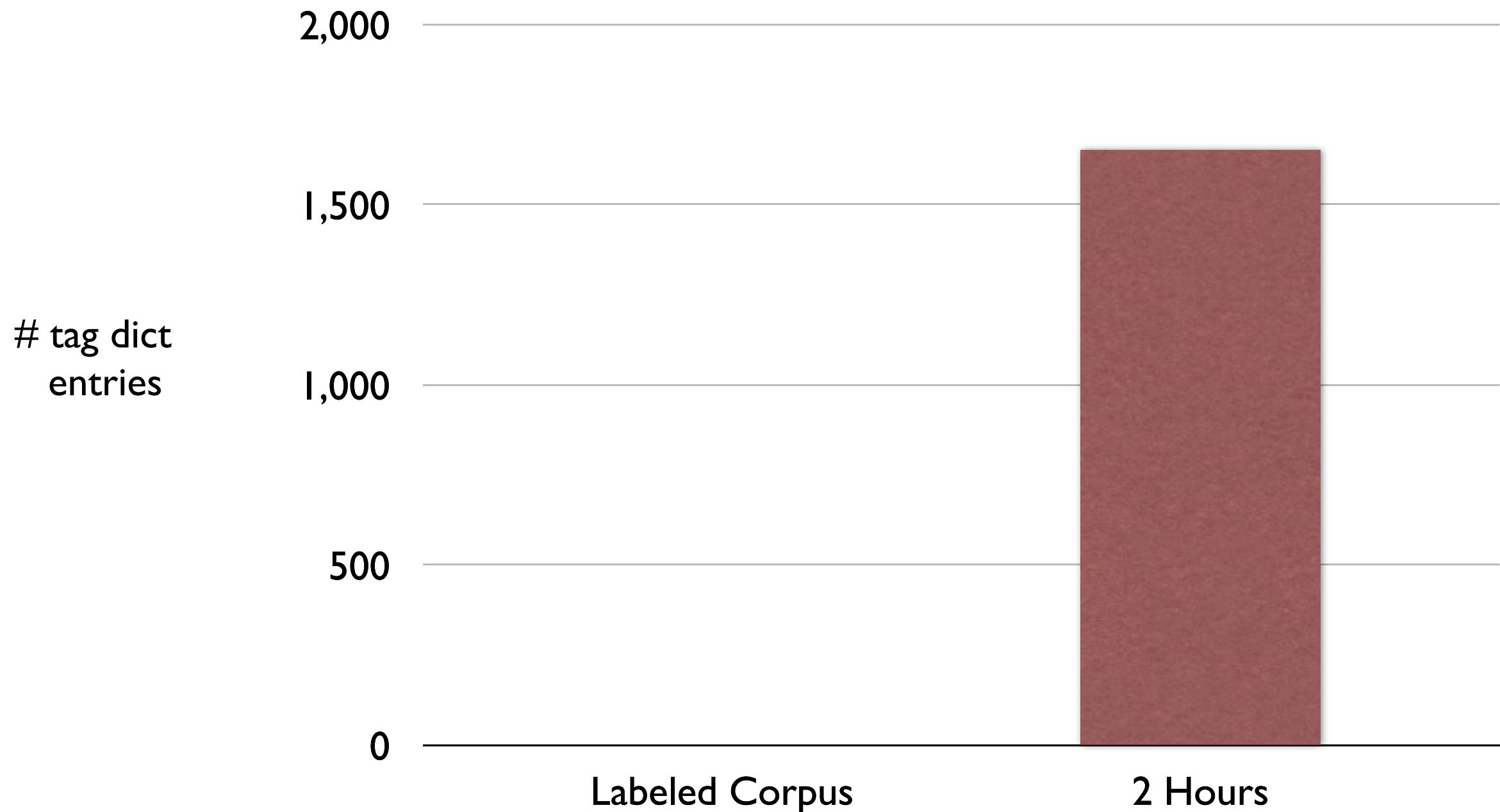
Tag dictionary

[Kupiec, 1992]
[Merialdo, 1994]

# 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]

# A **Real** Tag Dictionary

# tag dict
entries

# A **Real** Tag Dictionary

# A **Real** Tag Dictionary

# A **Real** Tag Dictionary

# A **Real** Tag Dictionary

Extremely low coverage means **most** words are **unknown**

# A **Real** Tag Dictionary

Extremely low coverage means **most** words are **unknown**

⇒ **Bad for learning**  (poorly constrained)

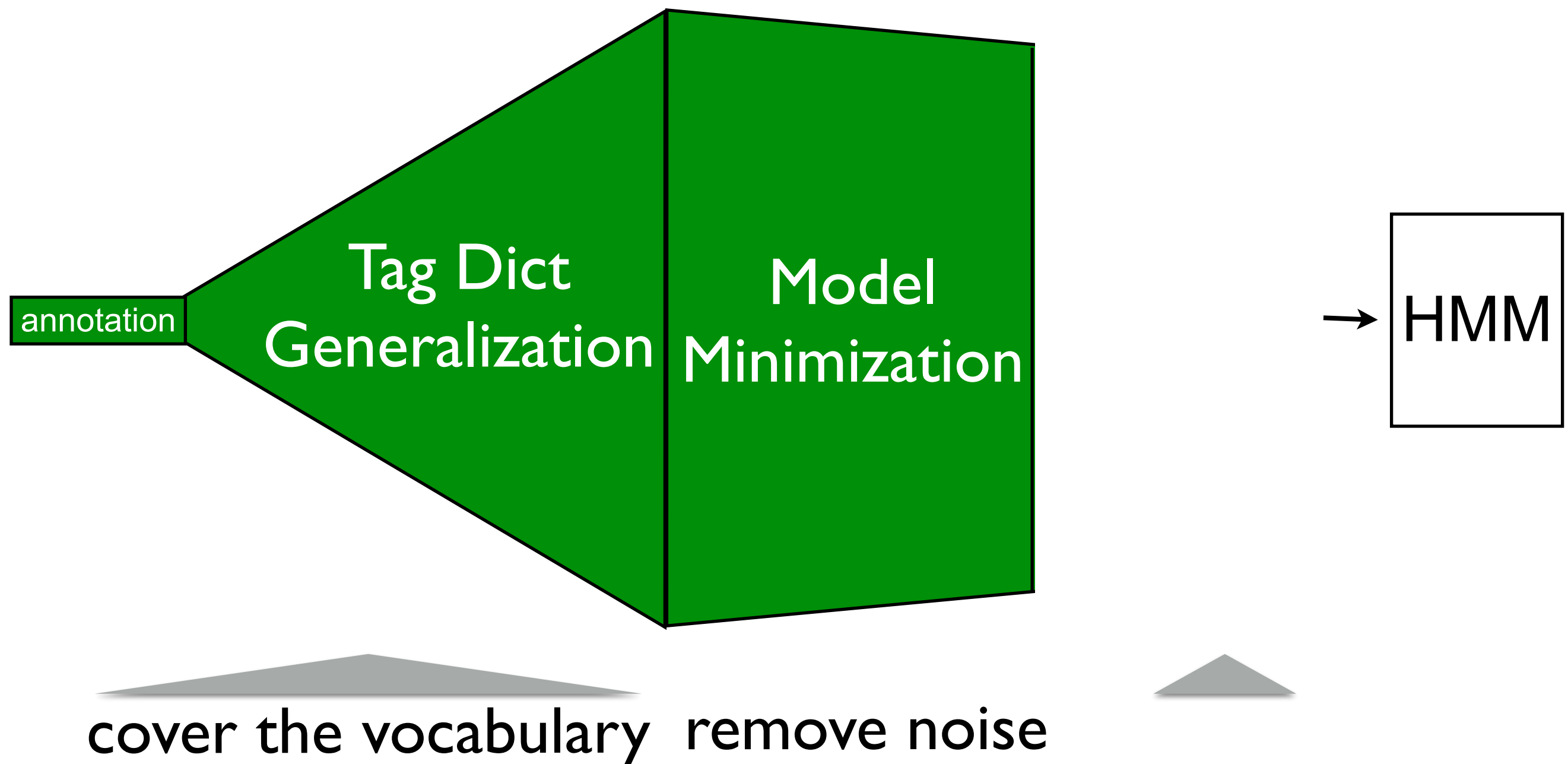# Our Approach
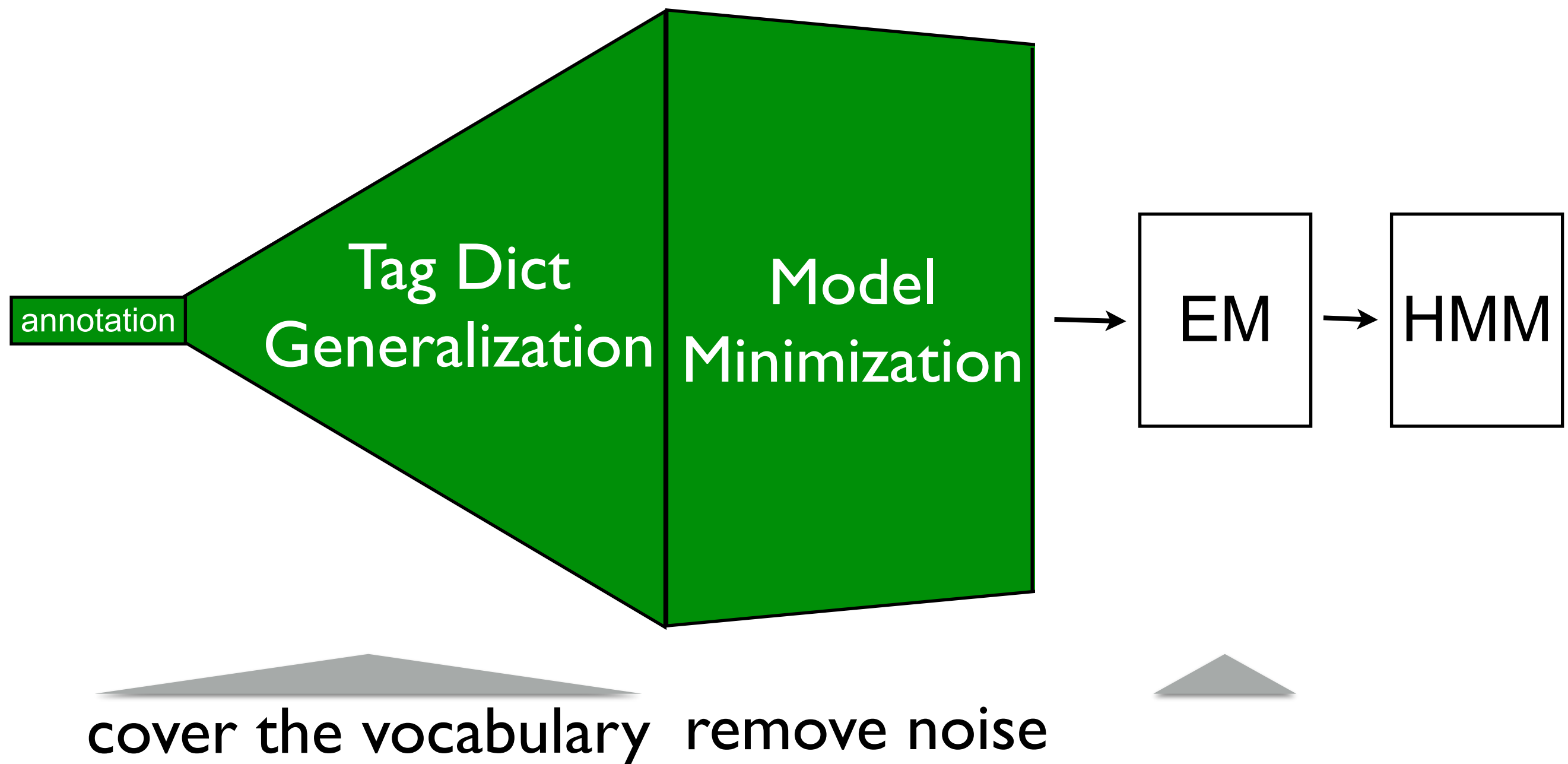
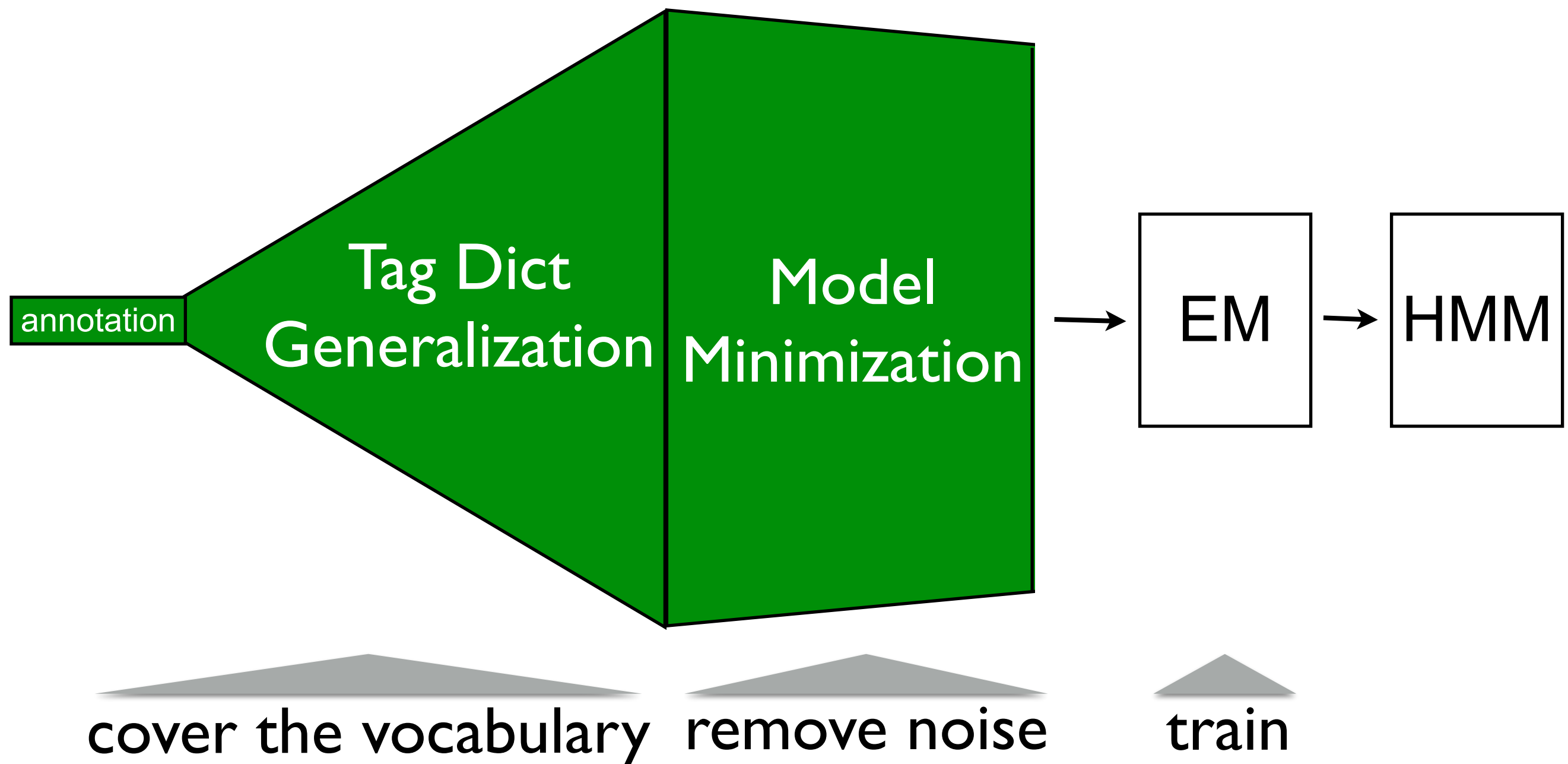annotation → HMM

# Our Approach

annotation → HMM

# Our Approach

# Our Approach

annotation

Tag Dict Generalization

→ HMM

cover the vocabulary

# Our Approach



annotation

Tag Dict Generalization

Model Minimization

→ HMM

cover the vocabulary  remove noise

# Our Approach



annotation
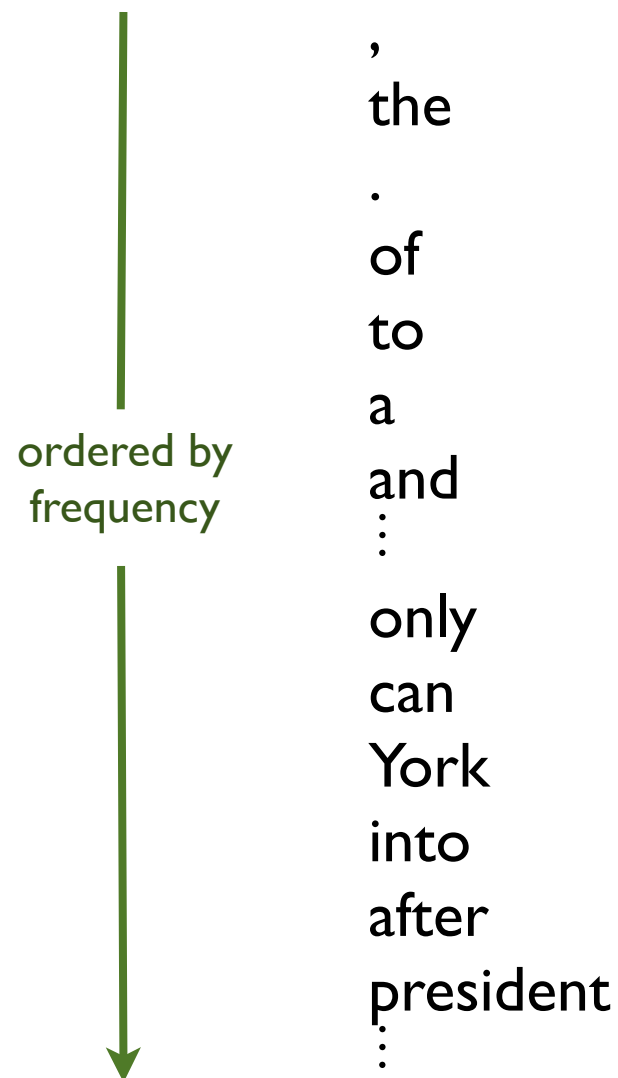
Tag Dict Generalization

Model Minimization

EM

HMM

cover the vocabulary    remove noise

# Our Approach

# Our Approach

# Our Approach

# Collecting Annotations

Task #1

**Up to 4 hours** to create a **tag dictionary**

# Collecting Annotations

Task #1

**Up to 4 hours** to create a **tag dictionary**

ordered by frequency

,
the
.
of
to
a
and
⋮
only
can
York
into
after
president
⋮

# Collecting Annotations

Task #1

**Up to 4 hours** to create a **tag dictionary**

| | | | |
|---|---|---|---|
| , | , | | |
| 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

**Up to 4 hours** to annotate **full sentences**

# Collecting Annotations

Task #2

**Up to 4 hours** to annotate **full sentences**

Pierre  Vinken  ,  61    years  old  ,  will  join  the  board  as  a    nonexecutive  director  Nov.    29    .

Mr.  Vinken      is    chairman  of    Elsevier    N.V. ,    the  Dutch    publishing  group      .

⋮

# Collecting Annotations

Task #2

**Up to 4 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 .

⋮

# Our Approach

annotation

Tag Dict Generalization

Model Minimization

EM

HMM

cover the vocabulary

remove noise

train

# Our Approach

# Tag Dict Generalization

These annotations are too sparse!

# Tag Dict Generalization

These annotations are too sparse!

➡️ Generalize to the entire vocabulary

# Tag Dict Generalization

Our strategy:  Label Propagation

[Talukdar and Crammer. 2009]

# Tag Dict Generalization

Our strategy:  Label Propagation

- **Connect** annotations to raw corpus tokens

[Talukdar and Crammer. 2009]

# 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

# Tag Dict Generalization

**Type Annotations**
the **DT**
dog **NN**

**Raw Corpus**

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

# Tag Dict Generalization

Type Annotations
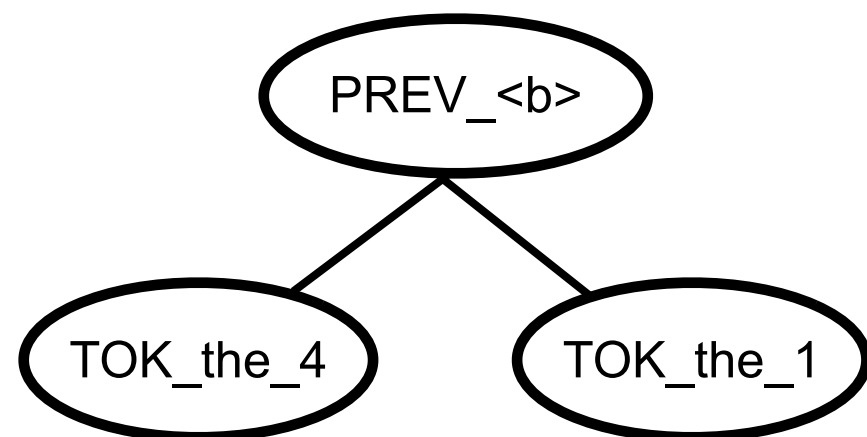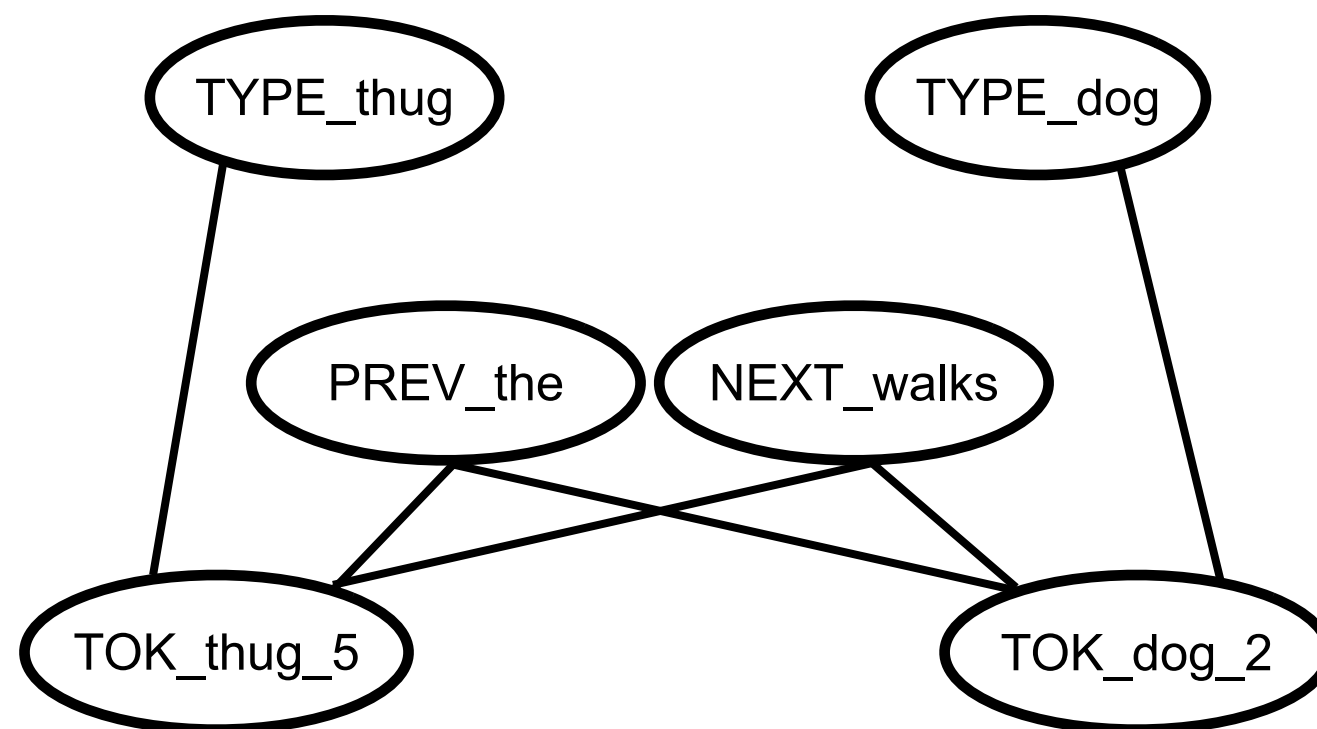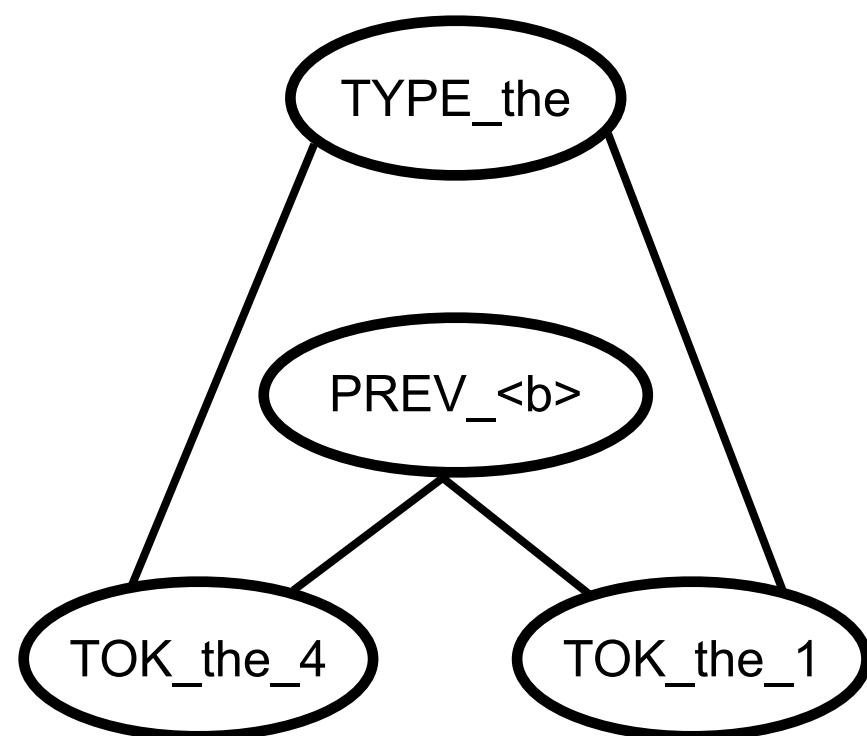the **DT**
dog **NN**

Raw Corpus
TOK_the_4

Token Annotations
th TOK_dog_2 ks
**DT** **NN** **VBZ**

# Tag Dict Generalization

Type Annotations
the   **DT**
dog  **NN**

Raw Corpus

TOK_the_4    TOK_the_1    TOK_thug_5    TOK_dog_2

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

# Tag Dict Generalization

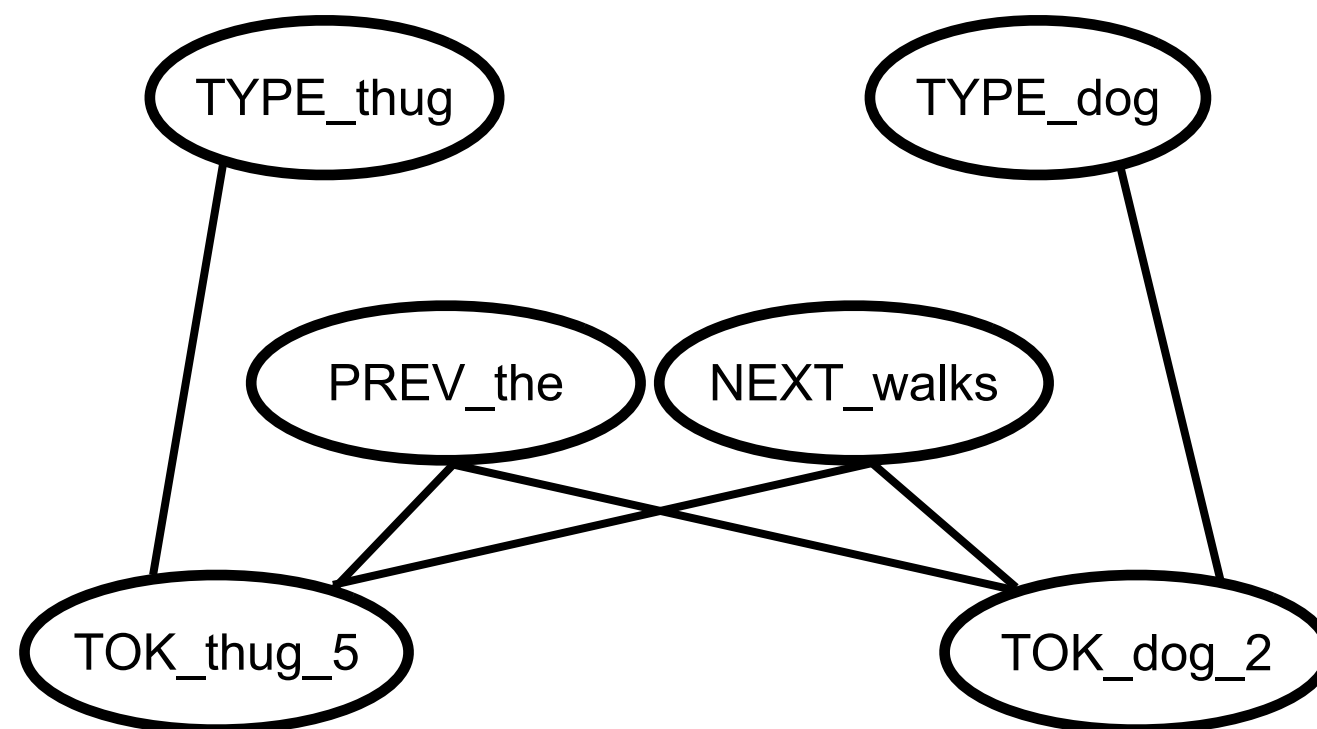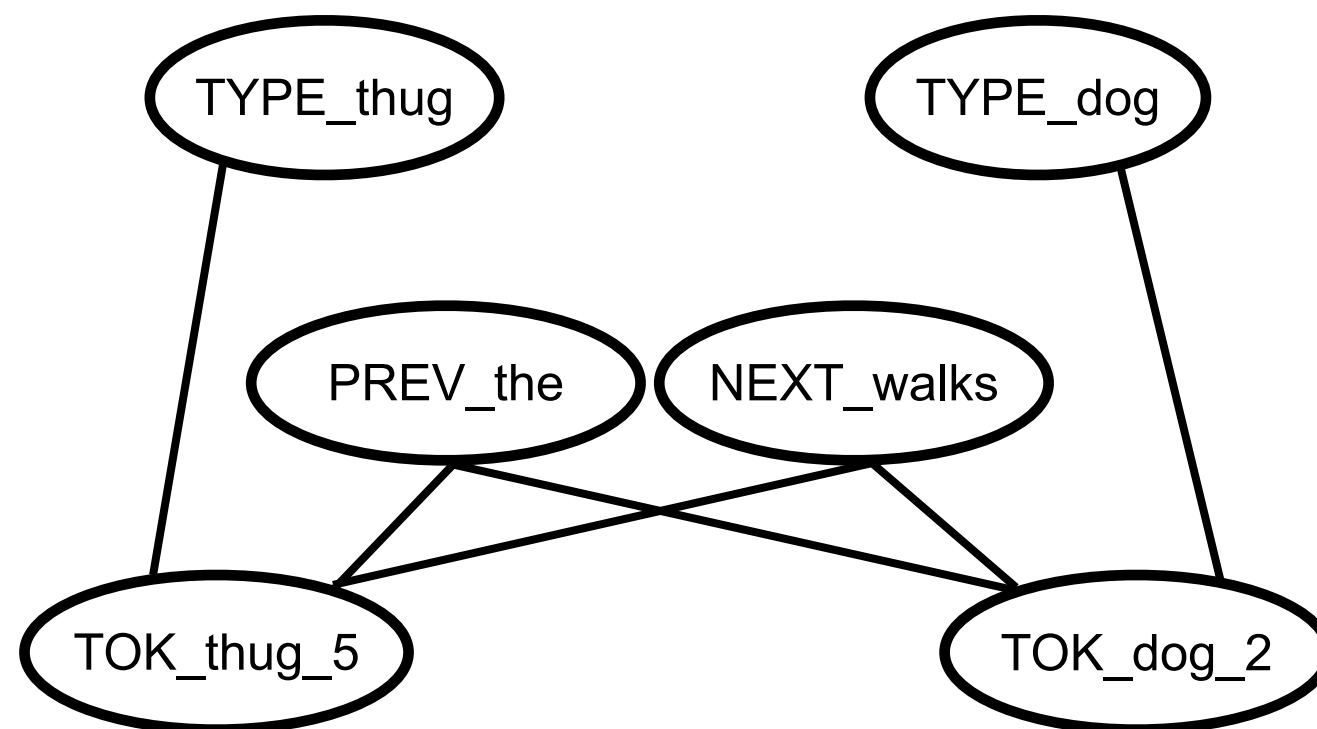Type Annotations
the   **DT**
dog  **NN**

Raw Corpus

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 **DT**
dog **NN**

Raw Corpus

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   **DT**
dog  **NN**

Raw Corpus

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 **DT**
dog **NN**

Raw Corpus

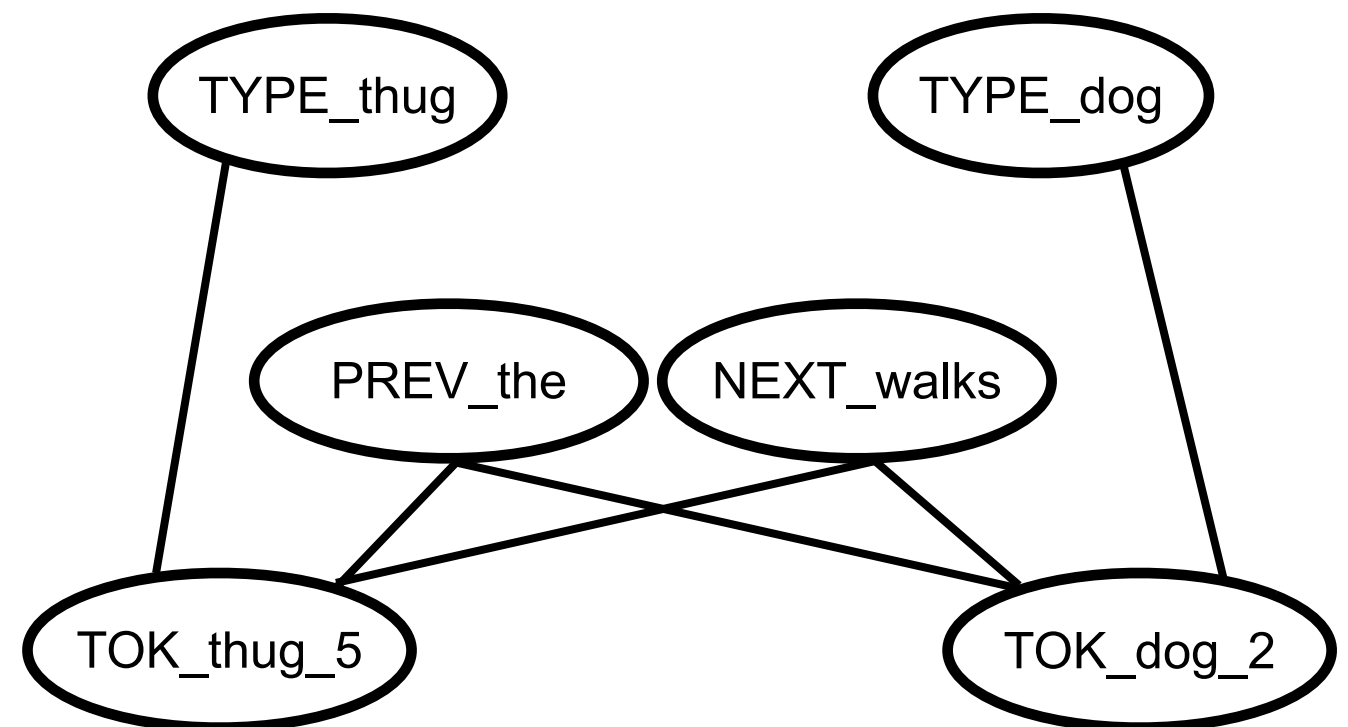TOK_the_4    TOK_the_1    TOK_thug_5    TOK_dog_2

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

# Tag Dict Generalization

Type Annotations
the **DT**
dog **NN**

Raw Corpus

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

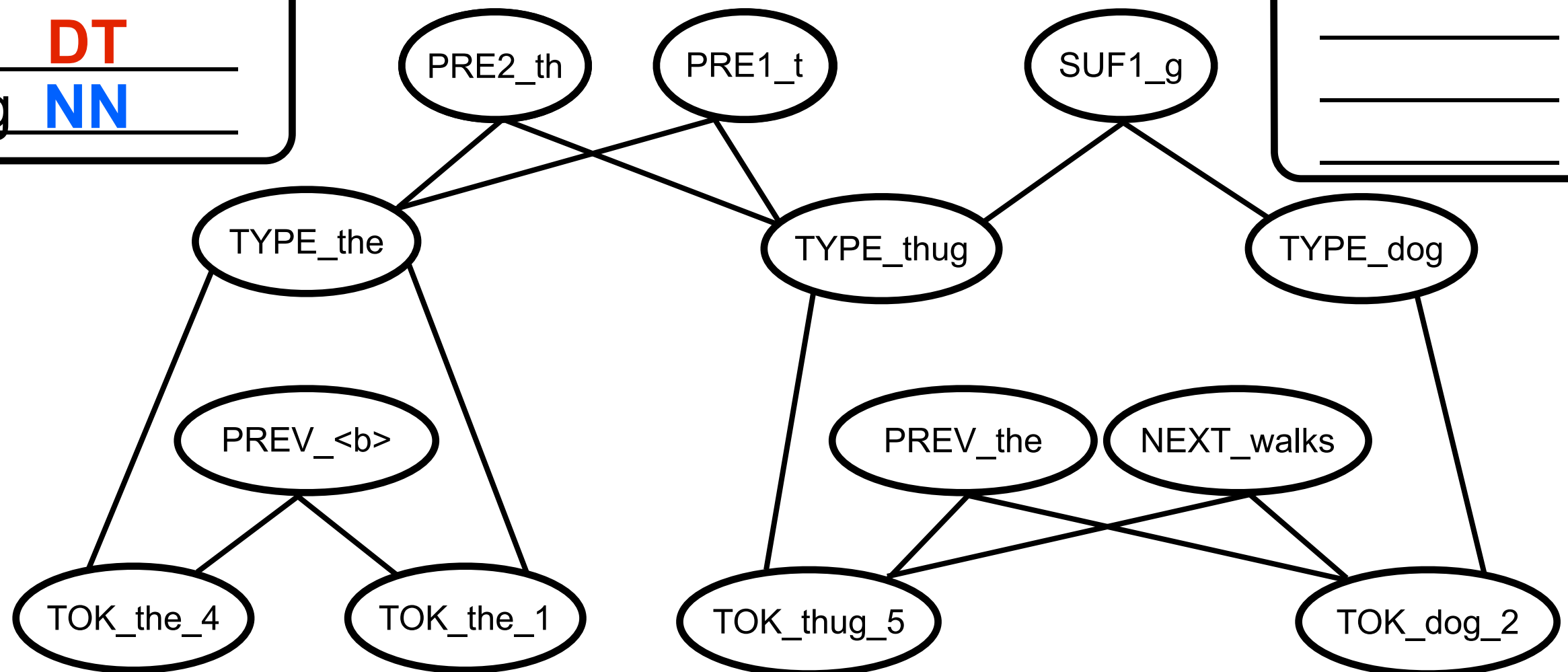Type Annotations
the **DT**
dog **NN**

Raw Corpus

PREV_\<b>

TOK_the_4    TOK_the_1

PREV_the    NEXT_walks

TOK_thug_5    TOK_dog_2

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

# Tag Dict Generalization

Type Annotations
the **DT**
dog **NN**

Raw Corpus

PREV_\<b\>

TOK_the_4   TOK_the_1

PREV_the   NEXT_walks

TOK_thug_5   TOK_dog_2

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

# Tag Dict Generalization

Type Annotations
th ( TYPE_dog )
dog **NN**

Raw Corpus

PREV_<b>

TOK_the_4        TOK_the_1

PREV_the        NEXT_walks

TOK_thug_5        TOK_dog_2

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

# Tag Dict Generalization

Type Annotations
the **DT**
dog **NN**

Raw Corpus

TYPE_the

TYPE_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 **DT**
dog **NN**

Raw Corpus

TYPE_the

TYPE_thug

TYPE_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

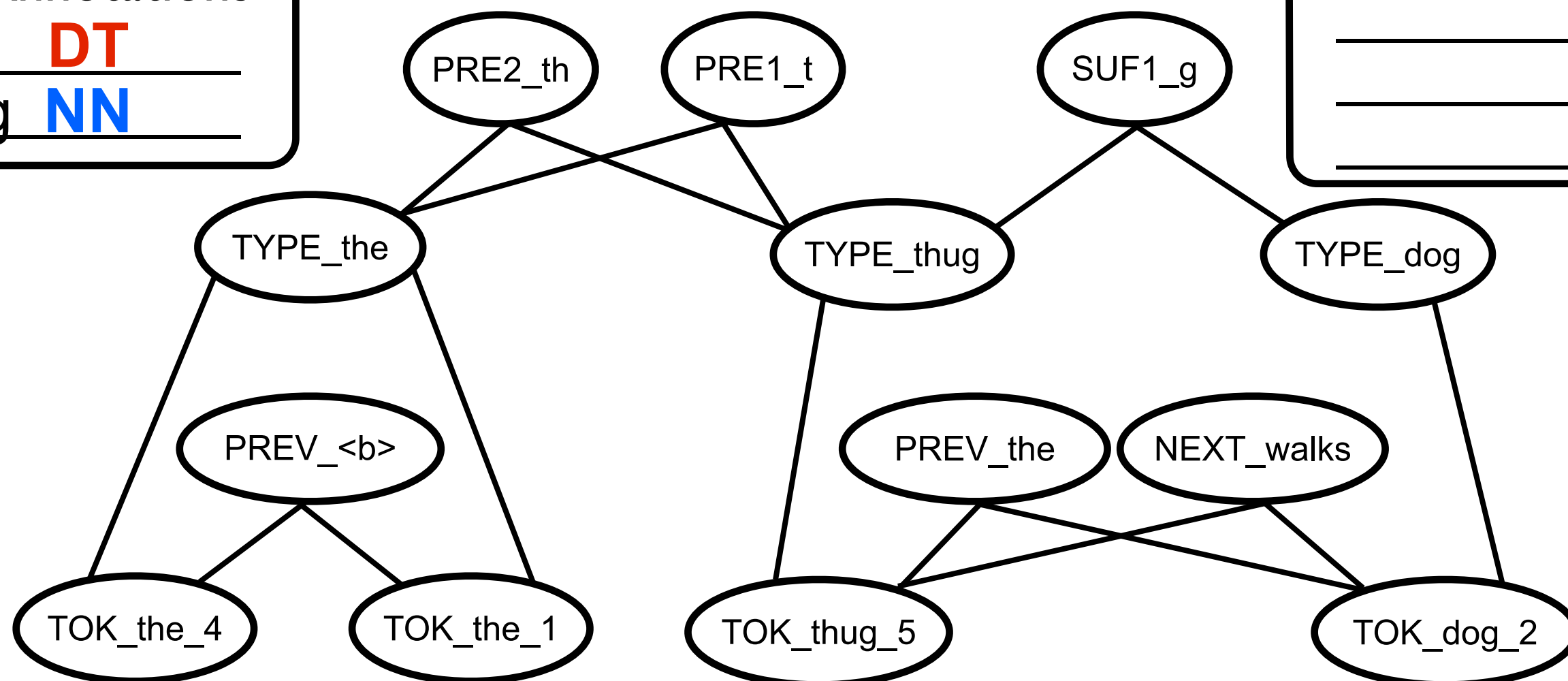# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization



**Type Annotations**
the   **DT**
dog  **NN**

**Raw Corpus**

TYPE_the

TYPE_thug

TYPE_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 **DT**
dog **NN**

PRE2_th

PRE1_t

SUF1_g

Raw Corpus

TYPE_the

TYPE_thug

TYPE_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

# Tag Dict Generalization

Type Annotations
the **DT**
dog **NN**

Raw Corpus

PRE2_th

PRE1_t

SUF1_g

TYPE_the

TYPE_thug

TYPE_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**

Any arbitrary features could be used

# Tag Dict Generalization

Type Annotations
the   **DT**
dog  **NN**

Raw Corpus

PRE2_th    PRE1_t    SUF1_g

TYPE_the    TYPE_thug    TYPE_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**

Any arbitrary features could be used

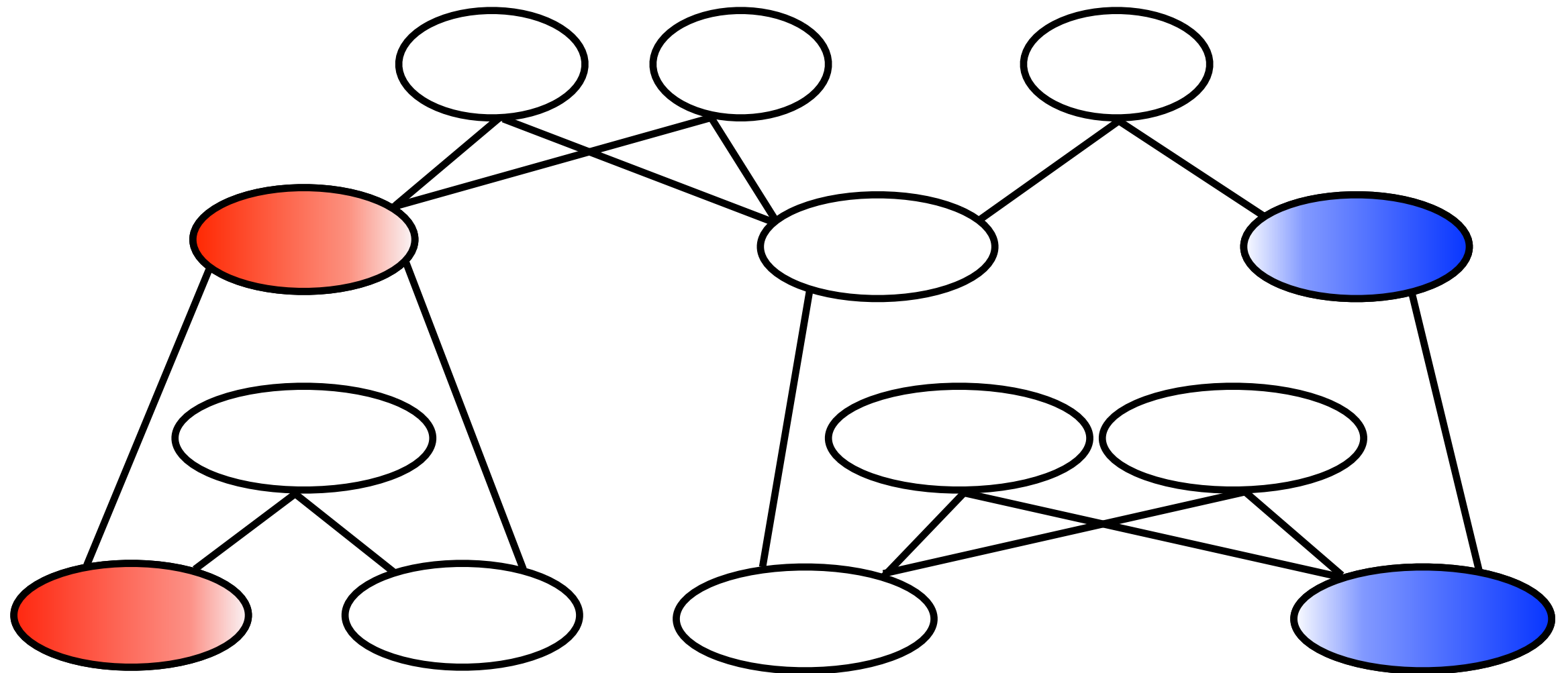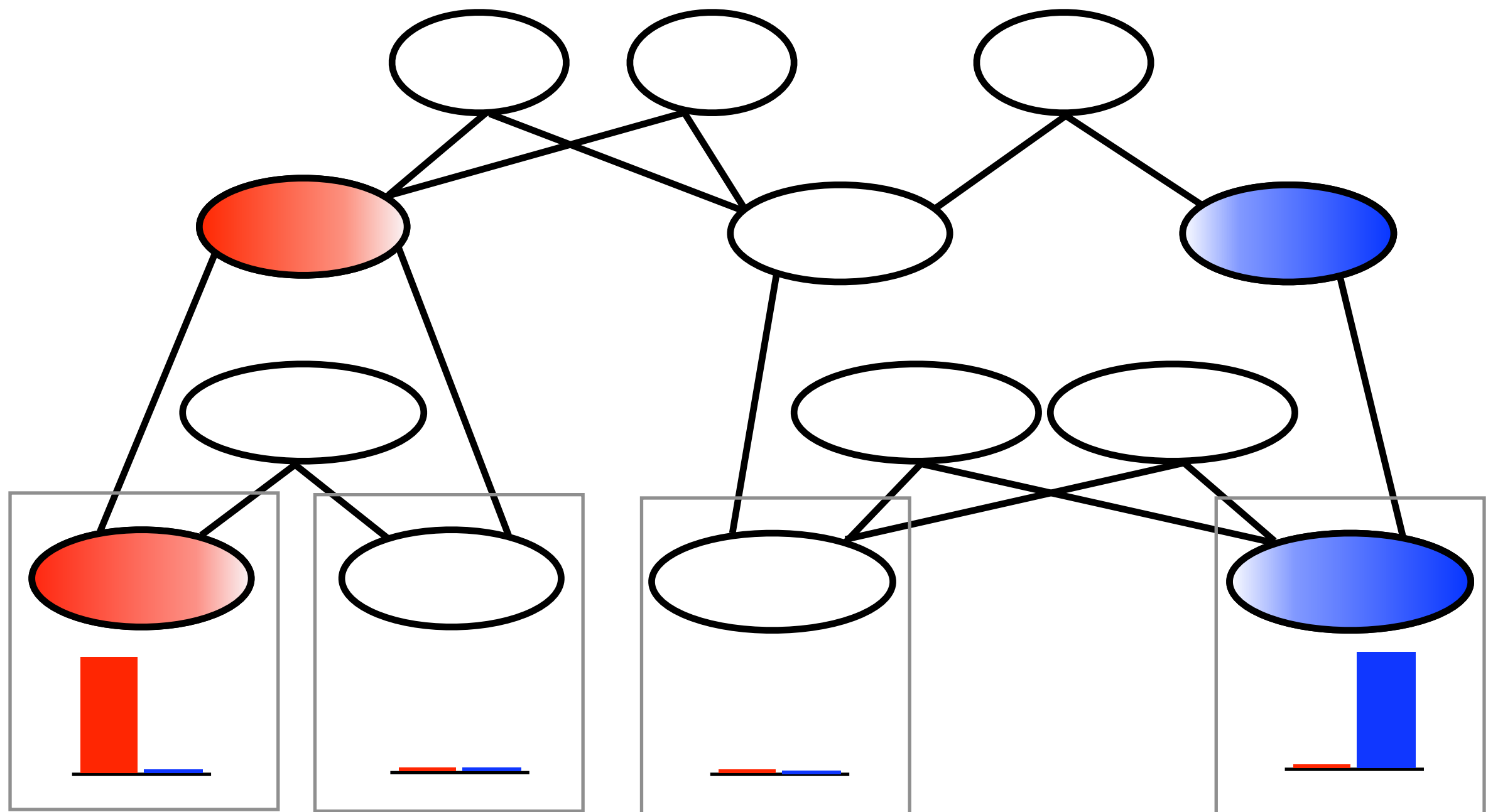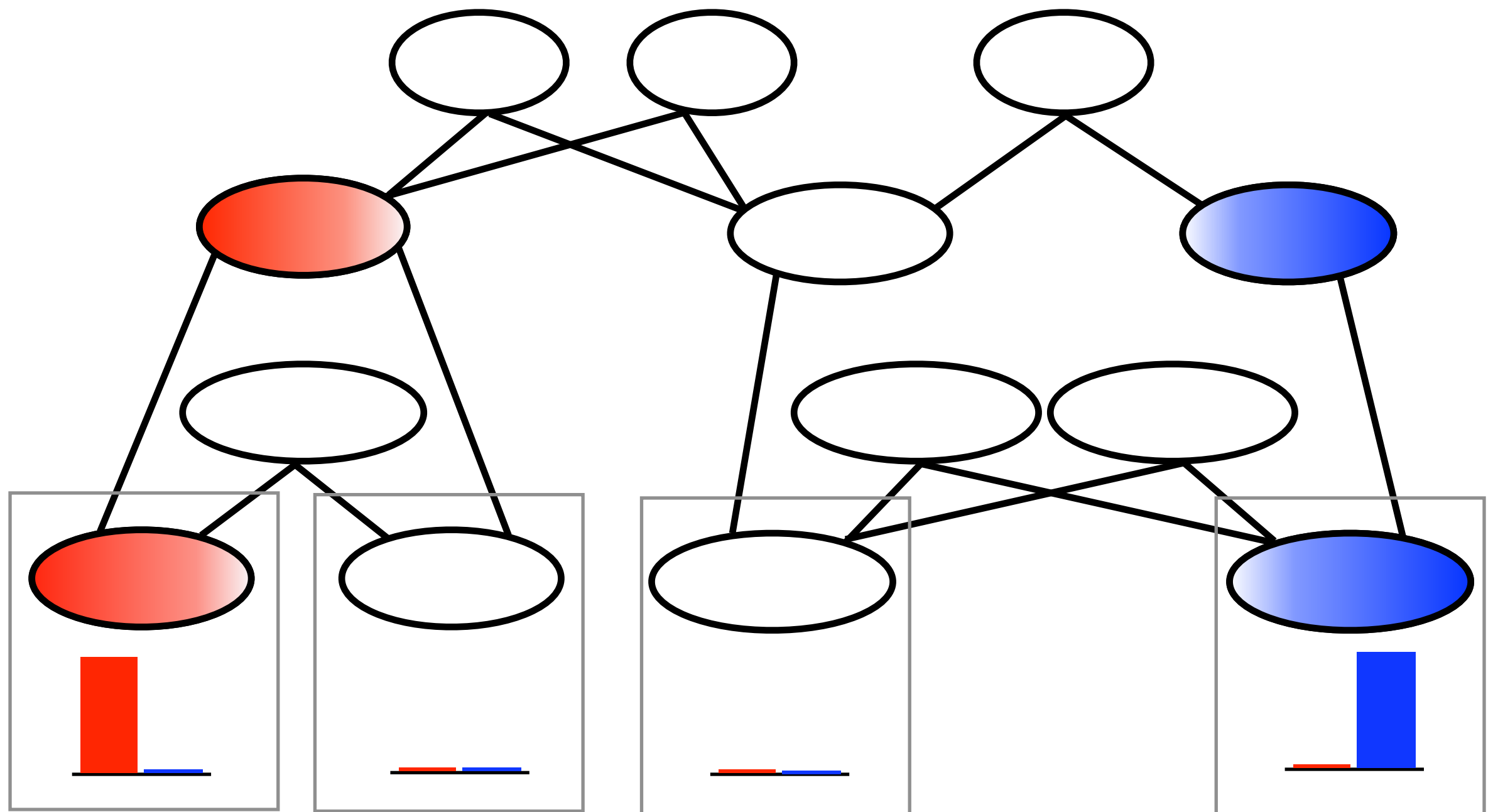# Tag Dict Generalization
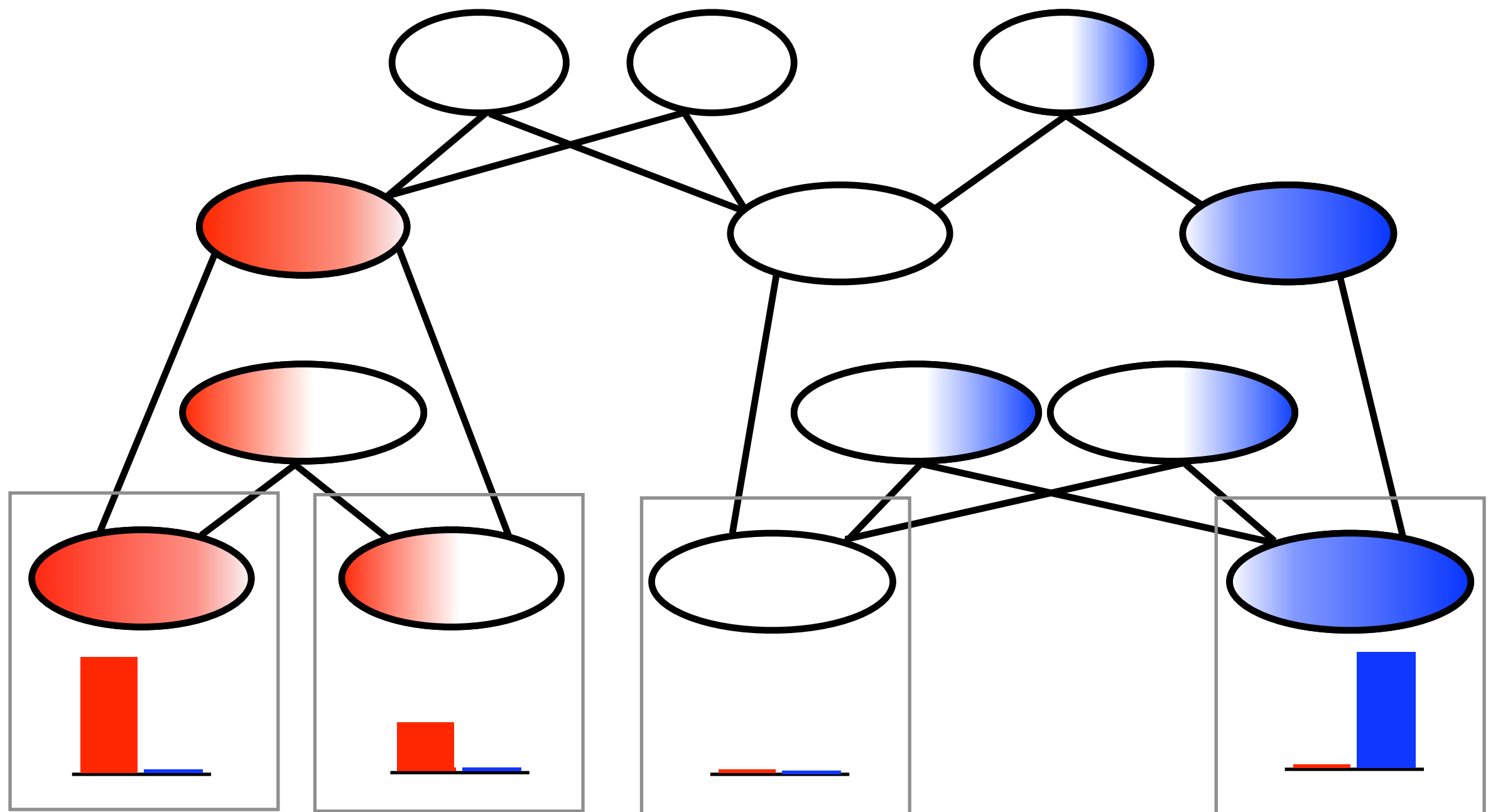
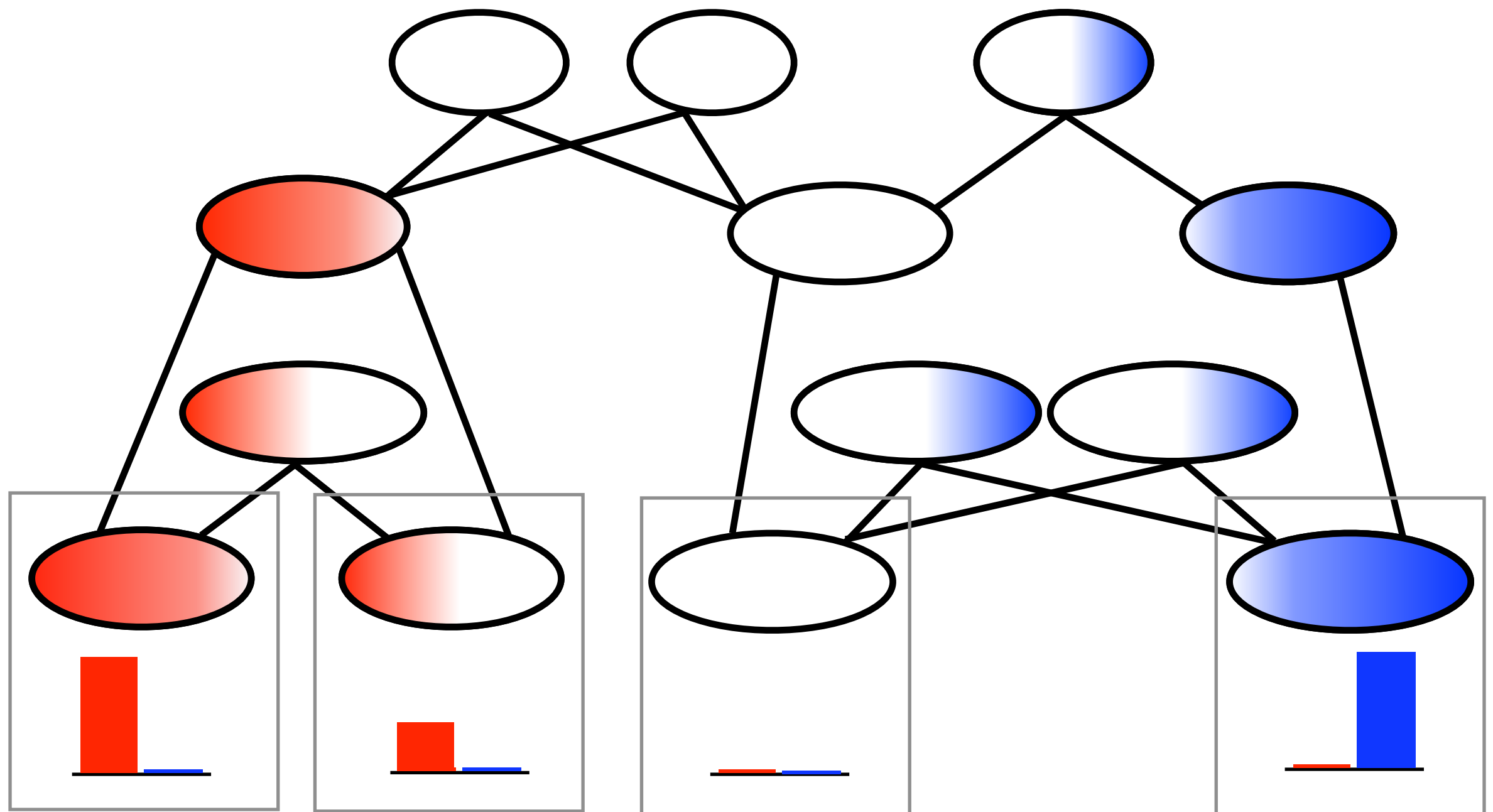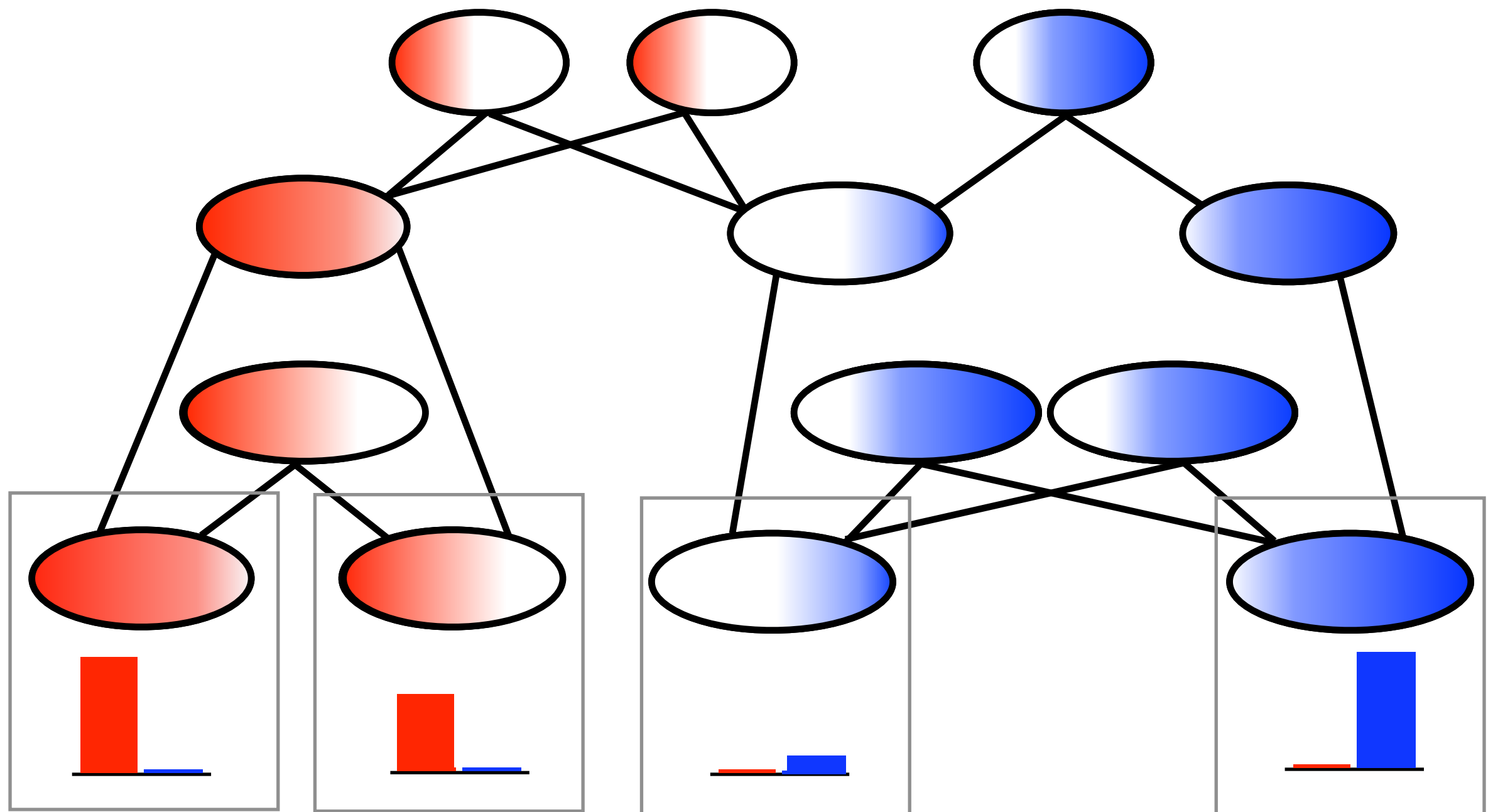# Tag Dict Generalization

# Tag Dict Generalization
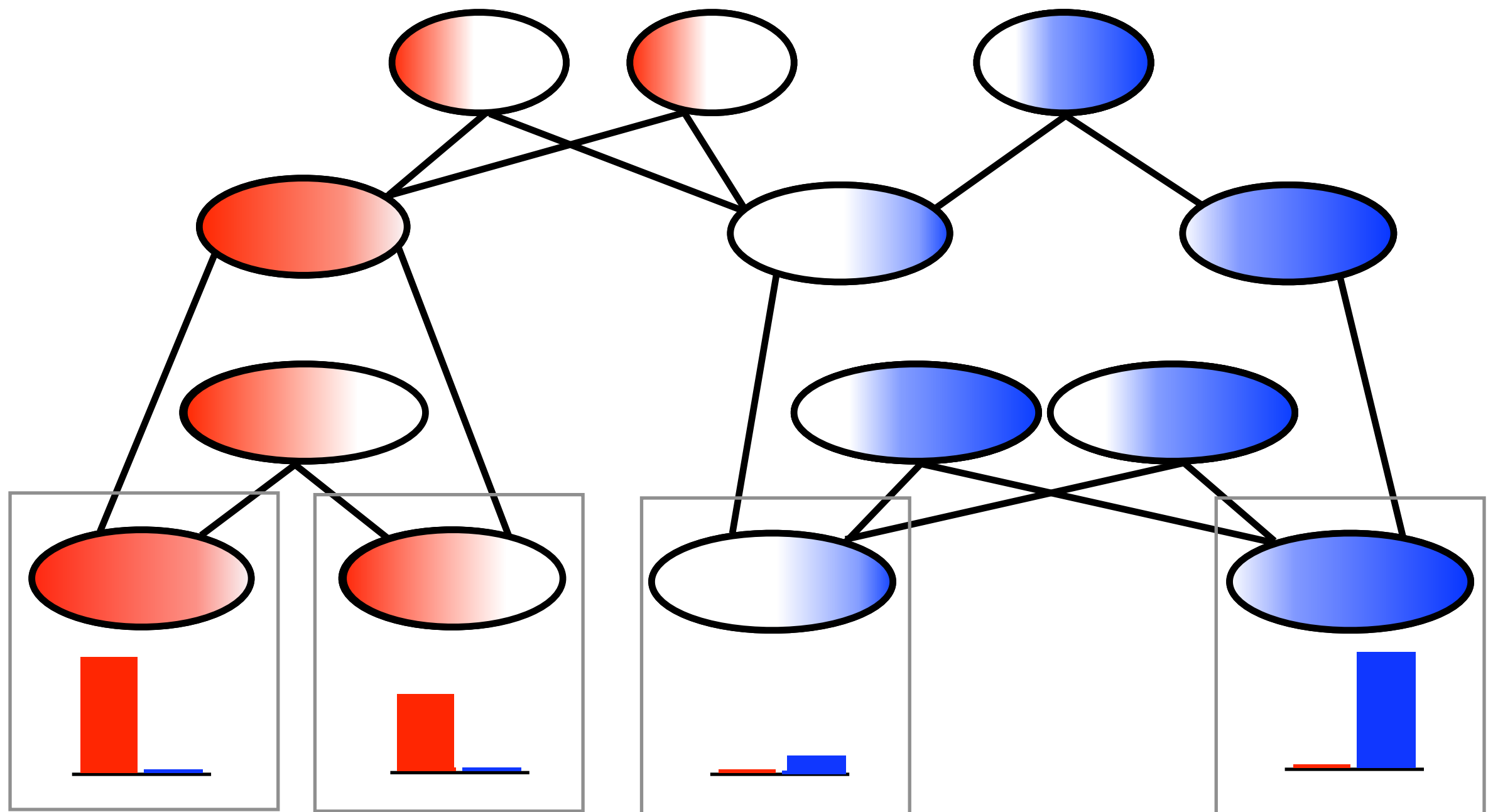
# Tag Dict Generalization

# Tag Dict Generalization

SUF1_g

TYPE_sibatarazuka

TYPE_dog

Finite-State Transducer (FST)

# Tag Dict Generalization



**Finite-State Transducer (FST)**

- Generates morphological analysis

# Tag Dict Generalization



Finite-State Transducer (FST)

- Generates morphological analysis

# Tag Dict Generalization



Finite-State Transducer (FST)

- Generates morphological analysis

# Tag Dict Generalization



Finite-State Transducer (FST)

- Generates morphological analysis

- Hand-built by a linguist in 10 hours

# Tag Dict Generalization



Finite-State Transducer (FST)

- Generates morphological analysis
- Hand-built by a linguist in 10 hours

# Tag Dict Generalization

SUF1_g

TYPE_dog

# Tag Dict Generalization

Type Annotations
the **DT**
dog **NN**

Raw Corpus

PRE2_th   PRE1_t   SUF1_g

TYPE_the   TYPE_thug   TYPE_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 **DT**

dog **NN**

Raw Corpus

PRE2_th    PRE1_t    SUF1_g

TYPE_the    TYPE_thug    TYPE_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 **DT**
dog **NN**

Raw Corpus

PRE2_th  PRE1_t  SUF1_g

TYPE_the  TYPE_thug  TYPE_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 **DT**
dog **NN**

Raw Corpus

PRE2_th    PRE1_t    SUF1_g

TYPE_the    TYPE_thug    TYPE_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 **DT**
dog **NN**

Raw Corpus

PRE2_th
PRE1_t
SUF1_g

TYPE_the
TYPE_thug
TYPE_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

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization
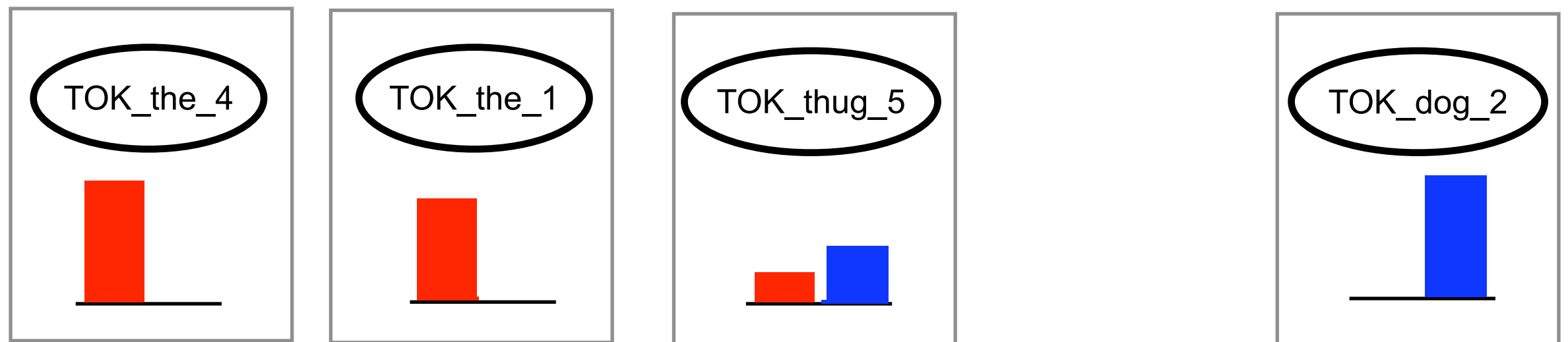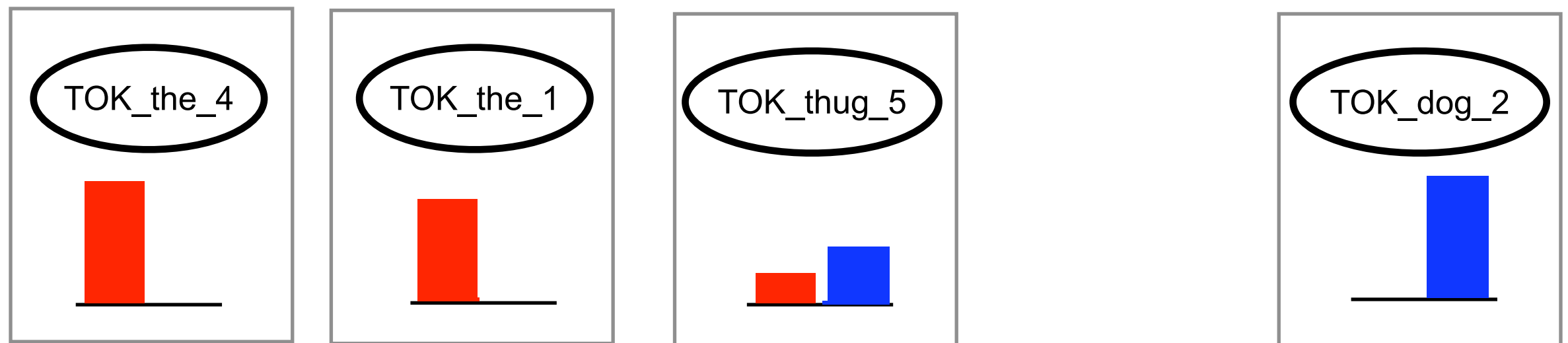
# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# Tag Dict Generalization

# 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

# Tag Dict Generalization

# Tag Dict Generalization

Result:

# Tag Dict Generalization

Result:

- a tag distribution on every token (soft tagging)

# 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.

- Approach based on work by Sujith Ravi and Kevin Knight (ISI)

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

# Model Minimization

# Model Minimization

$<b>_0$   The$_1$   man$_2$   saw$_3$   the$_4$   saw$_5$   $<b>_6$

# Model Minimization

$\text{<b>}_0$  $\text{The}_1$  $\text{man}_2$  $\text{saw}_3$  $\text{the}_4$  $\text{saw}_5$  $\text{<b>}_6$

<b>

**DT**

**NN**

**VBD**

# Model Minimization

DT NN NN DT NN
VBD VBD VBD

$\text{<b>}_0$ The$_1$ man$_2$ saw$_3$ the$_4$ saw$_5$ $\text{<b>}_6$

<b>

**DT**

**NN**

**VBD**

# Model Minimization

# Model Minimization

DT
NN
NN
DT
NN
VBD
VBD
VBD

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

$\langle b \rangle$
DT
NN
VBD

# Model Minimization

NN NN NN
DT VBD VBD DT VBD
$\langle b \rangle_0$ The$_1$ man$_2$ saw$_3$ the$_4$ saw$_5$ $\langle b \rangle_6$
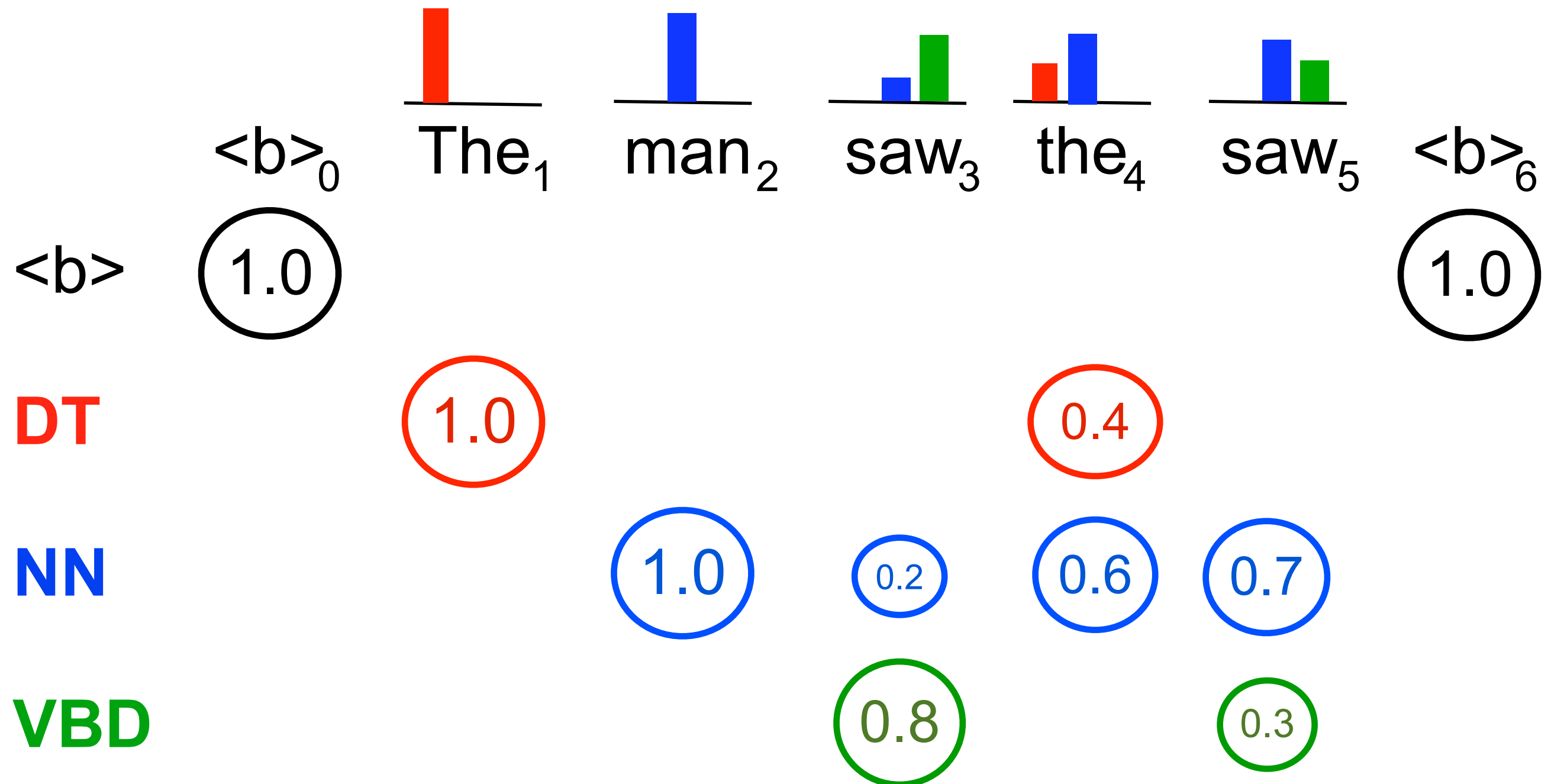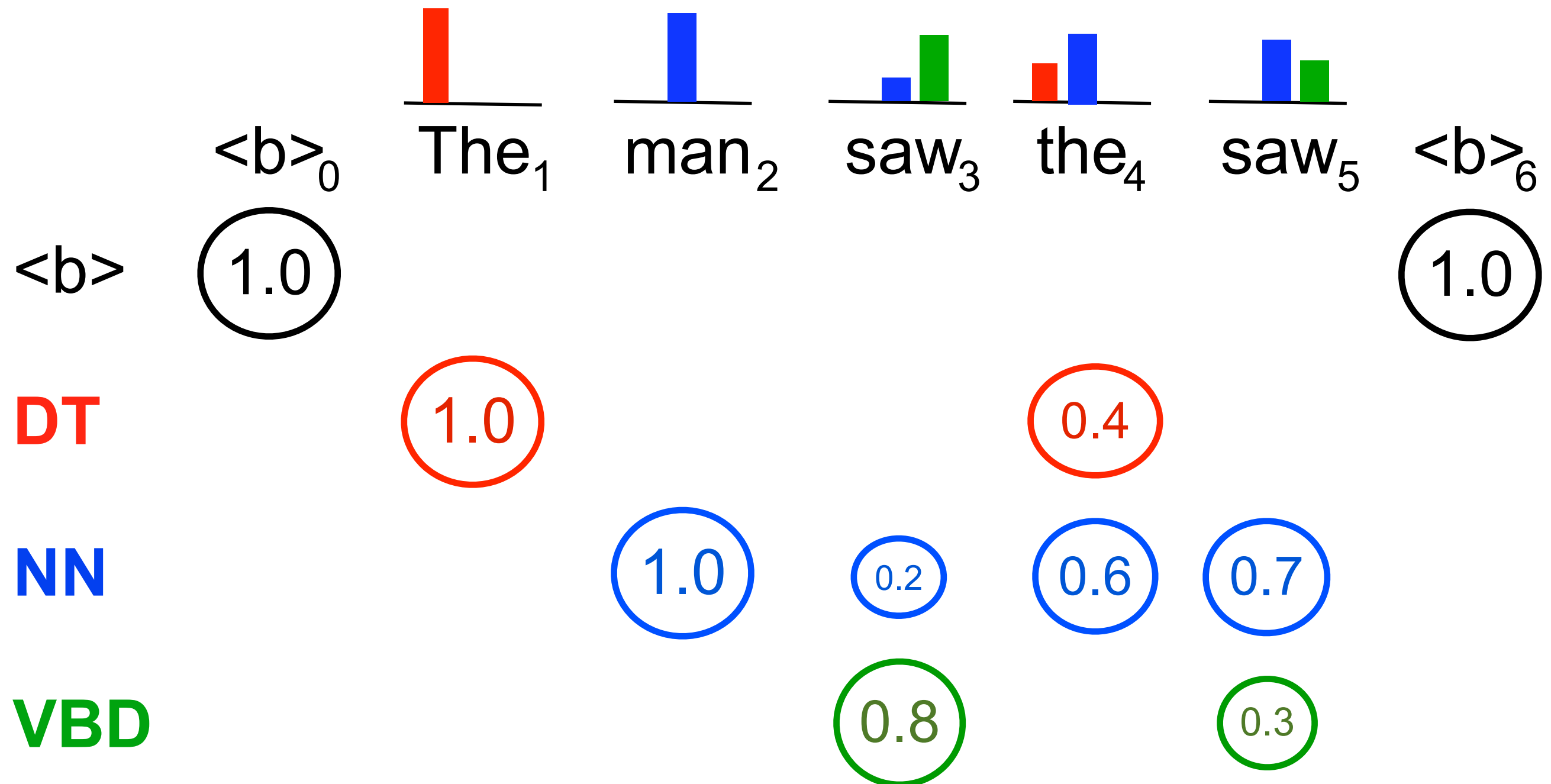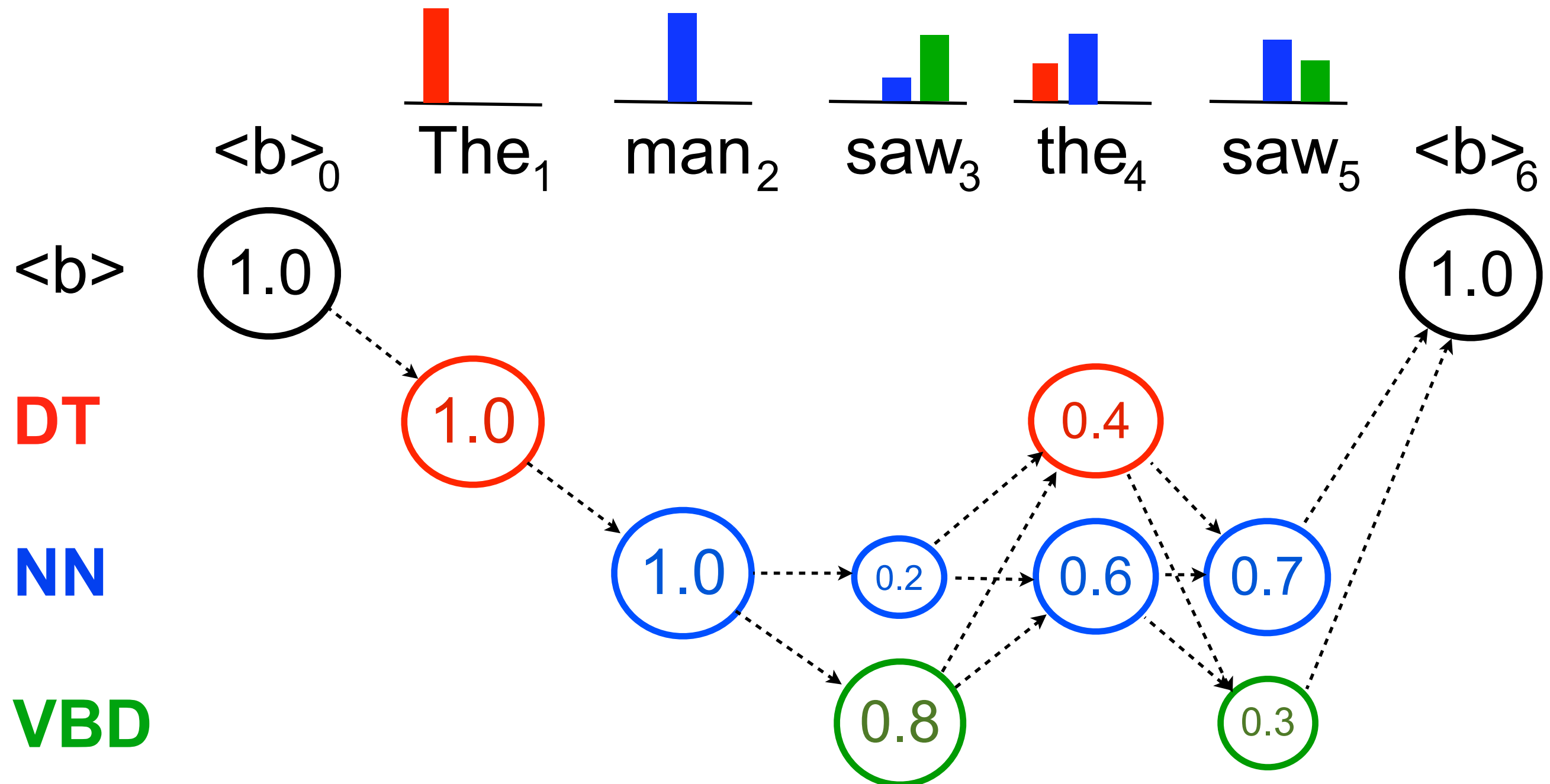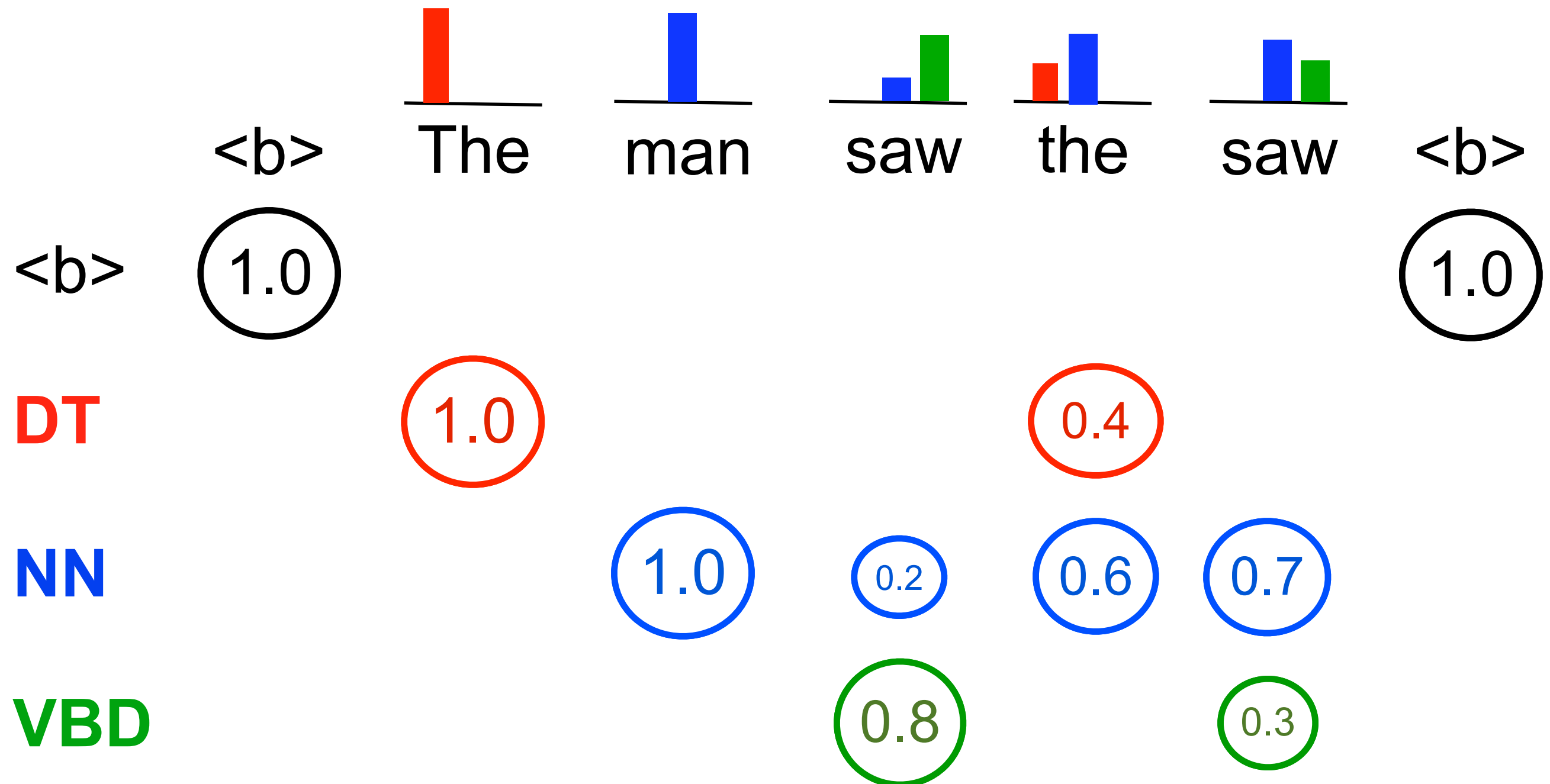
$\langle b \rangle$

DT

NN

VBD

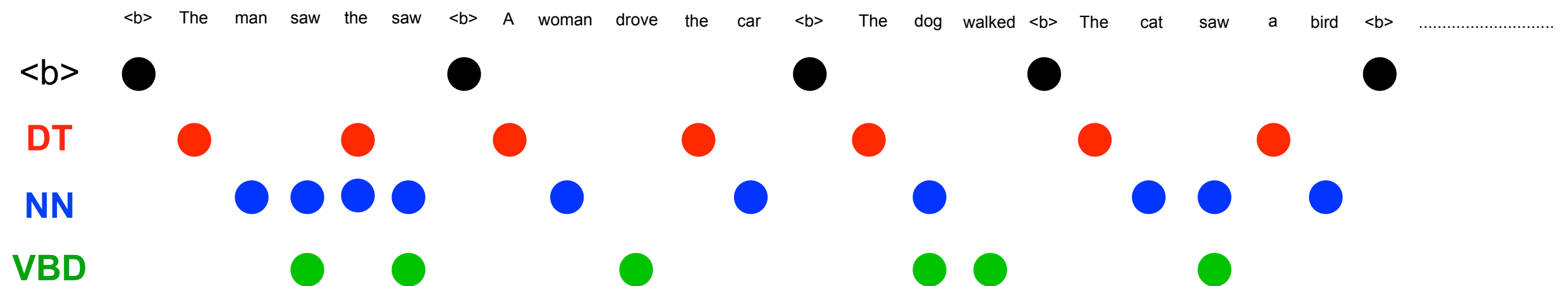# Model Minimization

# Model Minimization
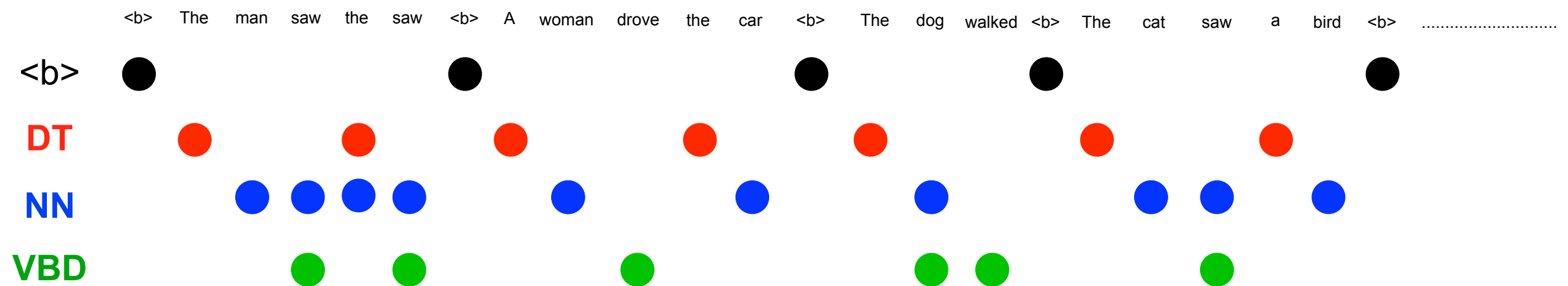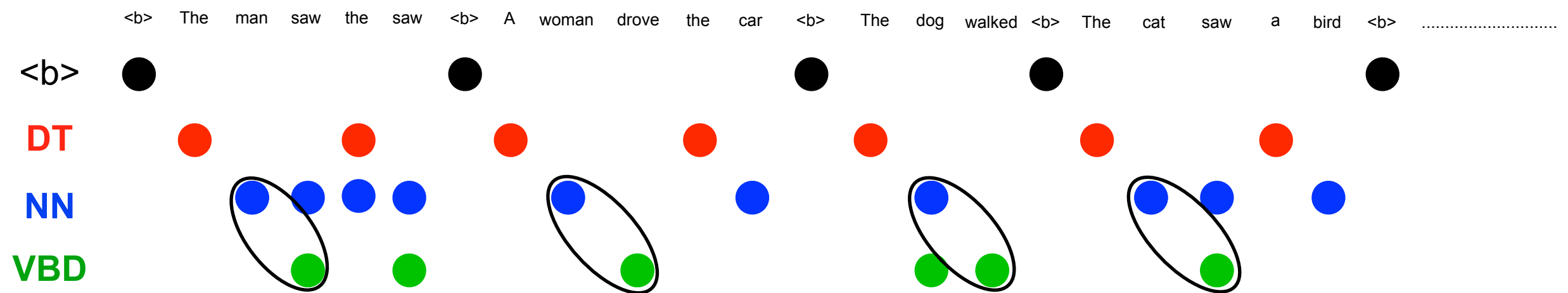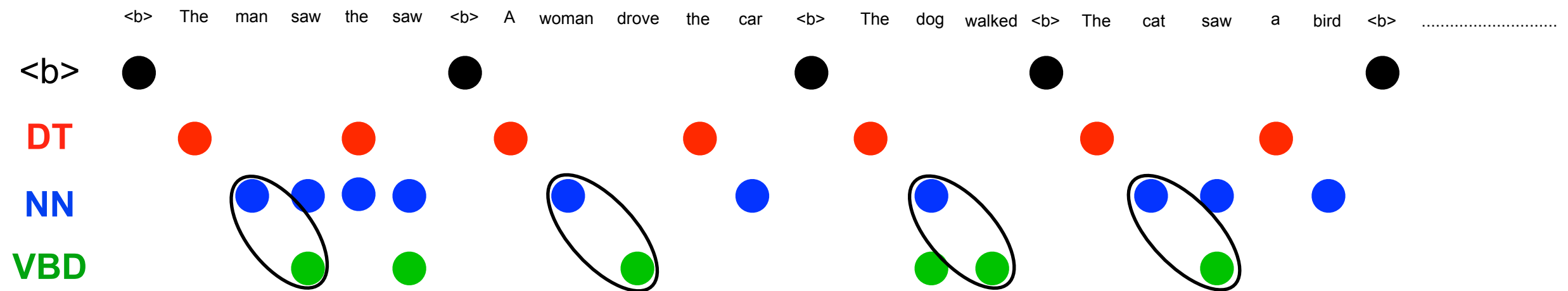
# Model Minimization

# Model Minimization

DT ? ? DT ?

<b>₀ The₁ man₂ saw₃ the₄ saw₅ <b>₆

<b>

DT

NN

VBD

# Model Minimization

**DT**     **?**     **?**     **DT**     **?**

$<b>_0$   The$_1$   man$_2$   saw$_3$   the$_4$   saw$_5$   $<b>_6$

<b>

**DT**

**NN**

**VBD**

# Model Minimization



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

<b>

**DT**

**NN**

**VBD**

# Model Minimization



&lt;b&gt;$_0$  The$_1$  man$_2$  saw$_3$  the$_4$  saw$_5$  &lt;b&gt;$_6$

&lt;b&gt;

**DT**

**NN**

**VBD**

# Model Minimization

# Model Minimization

Model Minimization

Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

Model Minimization

# Model Minimization
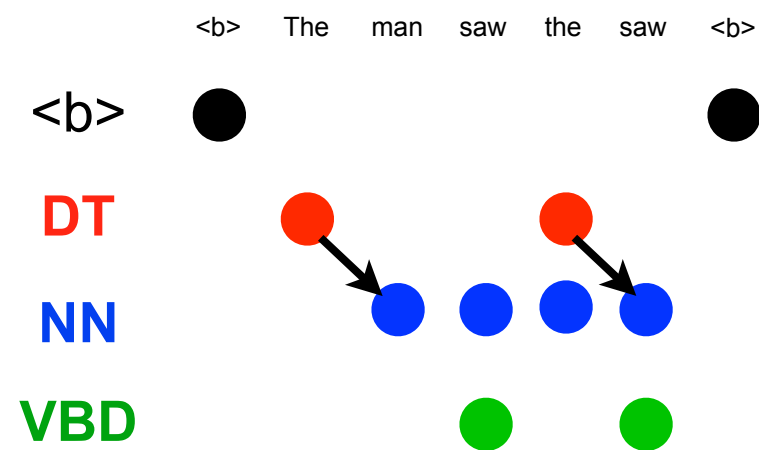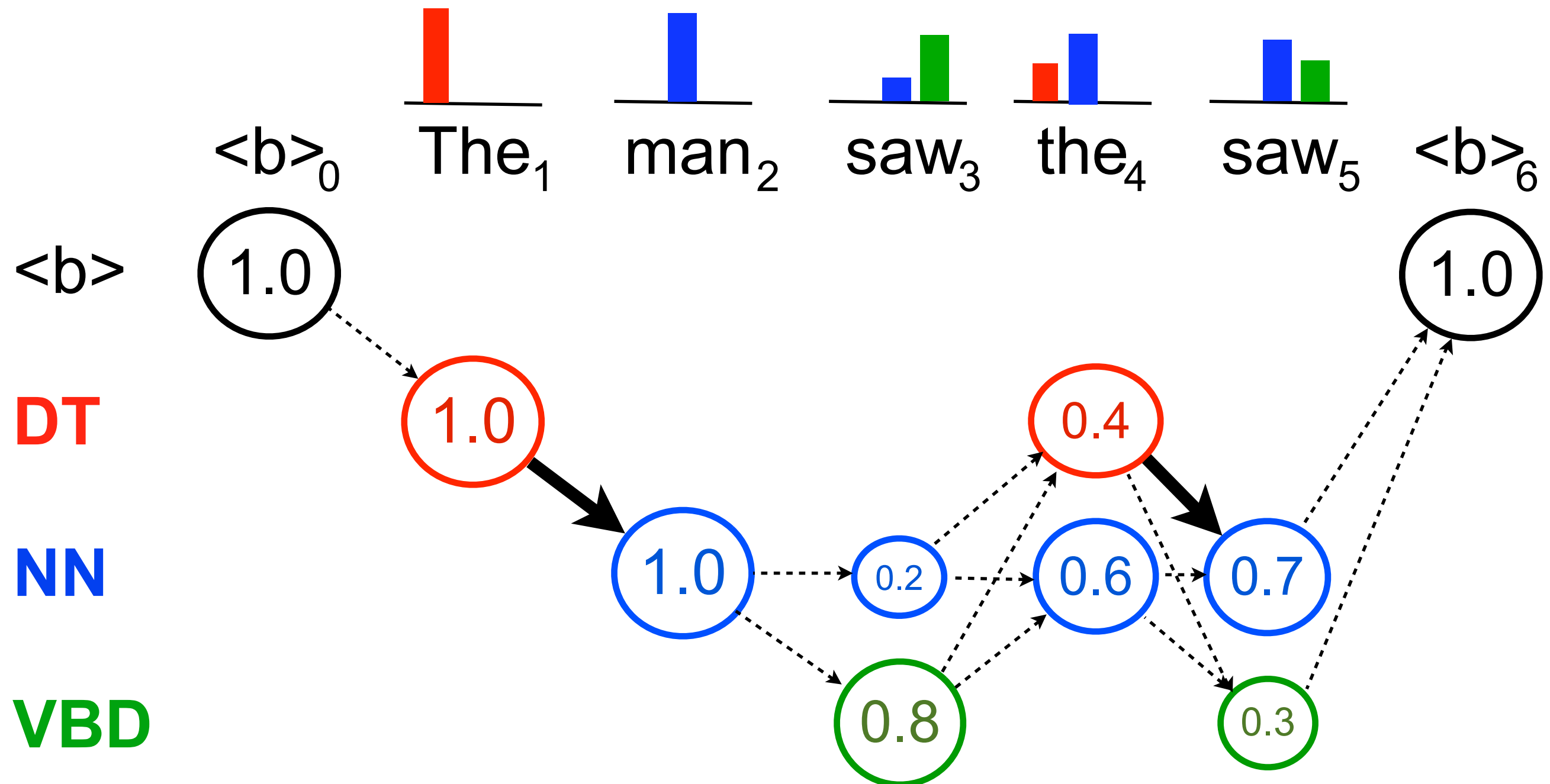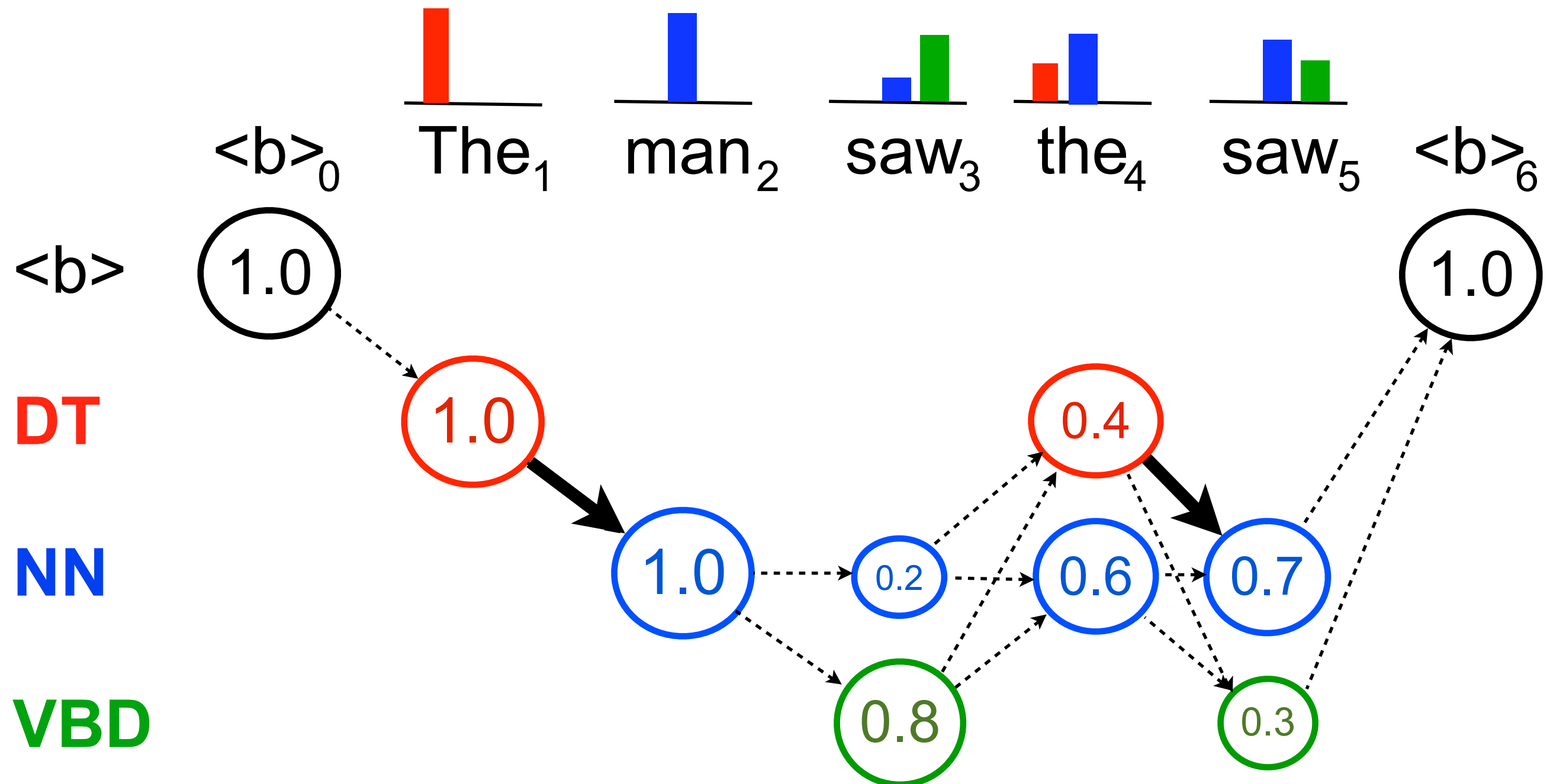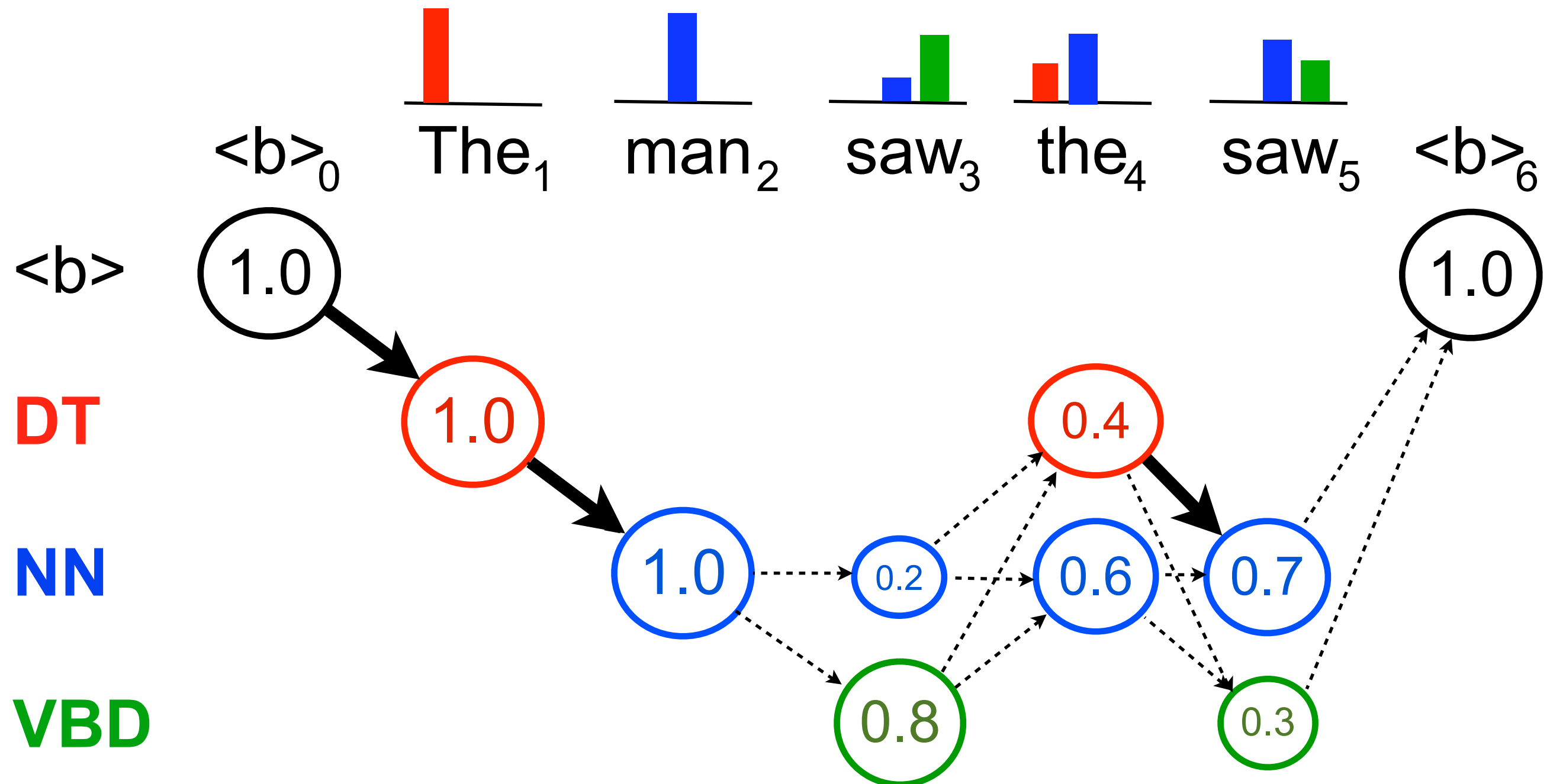
# Model Minimization

# Model Minimization

# Model Minimization

# Model Minimization

$<b>_0$ The$_1$ man$_2$ saw$_3$ the$_4$ saw$_5$ $<b>_6$

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

# Our Approach

# Our Approach

# EM Training

$\text{<b>}_0$ The$_1$ man$_2$ saw$_3$ the$_4$ saw$_5$ $\text{<b>}_6$

<span style="color:red">**DT**</span> <span style="color:blue">**NN**</span> <span style="color:green">**VBD**</span> <span style="color:red">**DT**</span> <span style="color:blue">**NN**</span>

# EM Training

Auto-Tagged
Corpus

# EM Training

Expanded
Tag Dictionary

_____

_____

_____

Auto-Tagged
Corpus

_____

_____

_____

# EM Training

Tag Dictionary
Generalization

Expanded
Tag Dictionary

_____

_____

_____

Auto-Tagged
Corpus

_____

_____

_____

# EM Training

Tag Dictionary Generalization → Expanded Tag Dictionary

Auto-Tagged Corpus → Initial Emissions

Auto-Tagged Corpus → Initial Transitions

# EM Training

# Results

# Types vs. Tokens

90 ─────────────────────────────

85 ─────────────────────────────

accuracy

80 ─────────────────────────────

75 ─────────────────────────────

70 ─────────────────────────────

[English]

# Types vs. Tokens



[English]

# Total Accuracy

| | English | Kinyarwanda | Malagasy |
|---|---|---|---|
| EM only | 69 | 69 | 77 |
| Our approach | 90 | 82 | 81 |

EM only   Our approach

[4 hours of type annotation]

# English Results

# English Results

All of **Wiktionary** (Li et al., 2012)

# English Results

All of **Wiktionary** (Li et al., 2012)          87%

# English Results

All of **Wiktionary** (Li et al., 2012)  87%

**Parallel Corpus** (Täckström et al., 2013)

# English Results

All of **Wiktionary** (Li et al., 2012)          87%

**Parallel Corpus** (Täckström et al., 2013)   89%

# English Results

All of **Wiktionary** (Li et al., 2012)          87%

**Parallel Corpus** (Täckström et al., 2013)   89%

**4-hours** (Garrette et al., 2013)

# English Results

All of **Wiktionary** (Li et al., 2012)       87%

**Parallel Corpus** (Täckström et al., 2013)  89%

**4-hours** (Garrette et al., 2013)           **90%**

# English Results

12
tags

All of **Wiktionary** (Li et al., 2012)          87%

**Parallel Corpus** (Täckström et al., 2013)  89%

**4-hours** (Garrette et al., 2013)          **90%**

# English Results

**12 tags**

All of **Wiktionary** (Li et al., 2012)    87%

**Parallel Corpus** (Täckström et al., 2013)    89%

**45 tags**

**4-hours** (Garrette et al., 2013)    **90%**

# Rich Morphology

# Rich Morphology

**Parallel Corpus** (Täckström et al., 2013)

Turkish

# Rich Morphology

**Parallel Corpus** (Täckström et al., 2013)

Turkish                                    65%

# Rich Morphology

**Parallel Corpus** (Täckström et al., 2013)

Turkish                    65%

**4-hours** (Garrette et al., 2013)

Kinyarwanda

# Rich Morphology

**Parallel Corpus** (Täckström et al., 2013)

Turkish                                            65%

**4-hours** (Garrette et al., 2013)

Kinyarwanda                             **82%**

# Current Work

- Minimally supervised CCG supertagging and parsing

- Human-provided GFL annotations

# Conclusion

- Our approach is able to achieve results better that or comparable to others, but given significantly less input.

- Our annotations are available to others.

- Software available as well.