# Feuding Families and Former Friends: Unsupervised Learning for Dynamic Fictional Relationships

Mohit Iyyer, Anupam Guha, Snigdha Chaturvedi, Jordan Boyd-Graber, and Hal Daumé III

University of Maryland, College Park
University of Colorado, Boulder

# How can we describe a fictional relationship between two characters?

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

Friend or foe?

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

<p align="center"><span style="color:blue">Friend</span> or <span style="color:red">foe</span>?</p>

<p align="center">Peter Pan and Captain Hook <em>(Peter Pan)</em></p>

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

<div align="center">

Friend or foe?

Peter Pan and Captain Hook *(Peter Pan)*

</div>

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

<div align="center">

Friend or foe?

Peter Pan and Captain Hook *(Peter Pan)*

Frodo and Sam (*Lord of the Rings*)

</div>

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

<div align="center">

**Friend** or **foe**?

Peter Pan and Captain Hook (*Peter Pan*)

Frodo and Sam (*Lord of the Rings*)

</div>

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

## Friend or foe?

Peter Pan and Captain Hook (*Peter Pan*)

Frodo and Sam (*Lord of the Rings*)

Winston and Julia (*1984*)

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

<p align="center">Friend or foe?</p>

<p align="center">Peter Pan and Captain Hook (<em>Peter Pan</em>)</p>

<p align="center">Frodo and Sam (<em>Lord of the Rings</em>)</p>

<p align="center">Winston and Julia (<em>1984</em>) <strong>???</strong></p>

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

Friend or foe?

Peter Pan and Captain Hook (*Peter Pan*)

Frodo and Sam (*Lord of the Rings*)

Winston and Julia (*1984*) **???**

Harry Potter and Sirius (*Prisoner of Azkaban*)

# How can we describe a fictional relationship between two characters?

- isn't this easy? we can assign it a single label (or *relationship descriptor*) from a predetermined set

<div align="center">

Friend or foe?

Peter Pan and Captain Hook (*Peter Pan*)

Frodo and Sam (*Lord of the Rings*)

Winston and Julia (*1984*) **???**

Harry Potter and Sirius (*Prisoner of Azkaban*) **???**

</div>

# How can we describe a fictional relationship between two characters?

- what if we treat relationships as sequences (or *trajectories*) of descriptors? (Chaturvedi et al., 2016)

Tom Sawyer and Becky Thatcher:
friends -> foes -> friends

# How can we describe a fictional relationship between two characters?

- what if we treat relationships as sequences (or *trajectories*) of descriptors? (Chaturvedi et al., 2016)
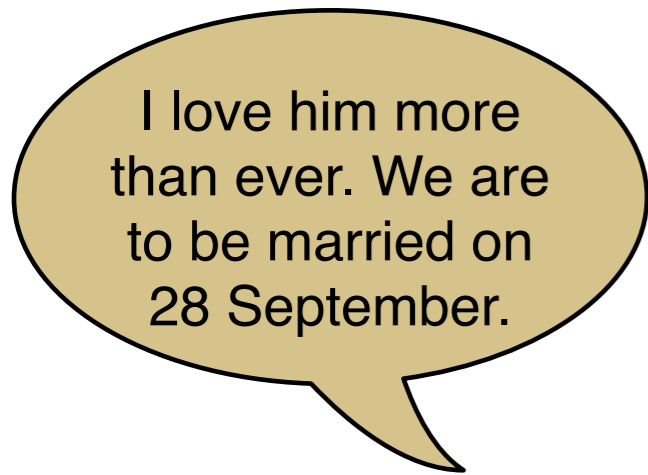
Tom Sawyer and Becky Thatcher:
friends -> foes -> friends

- limited by fixed descriptor set

- required expensive annotations

- limited to plot summaries

passage of time

Arthur and Lucy (*Dracula*)

I love him more than ever. We are to be married on 28 September.

passage of time

Arthur and Lucy (*Dracula*)

joy

I love him more than ever. We are to be married on 28 September.
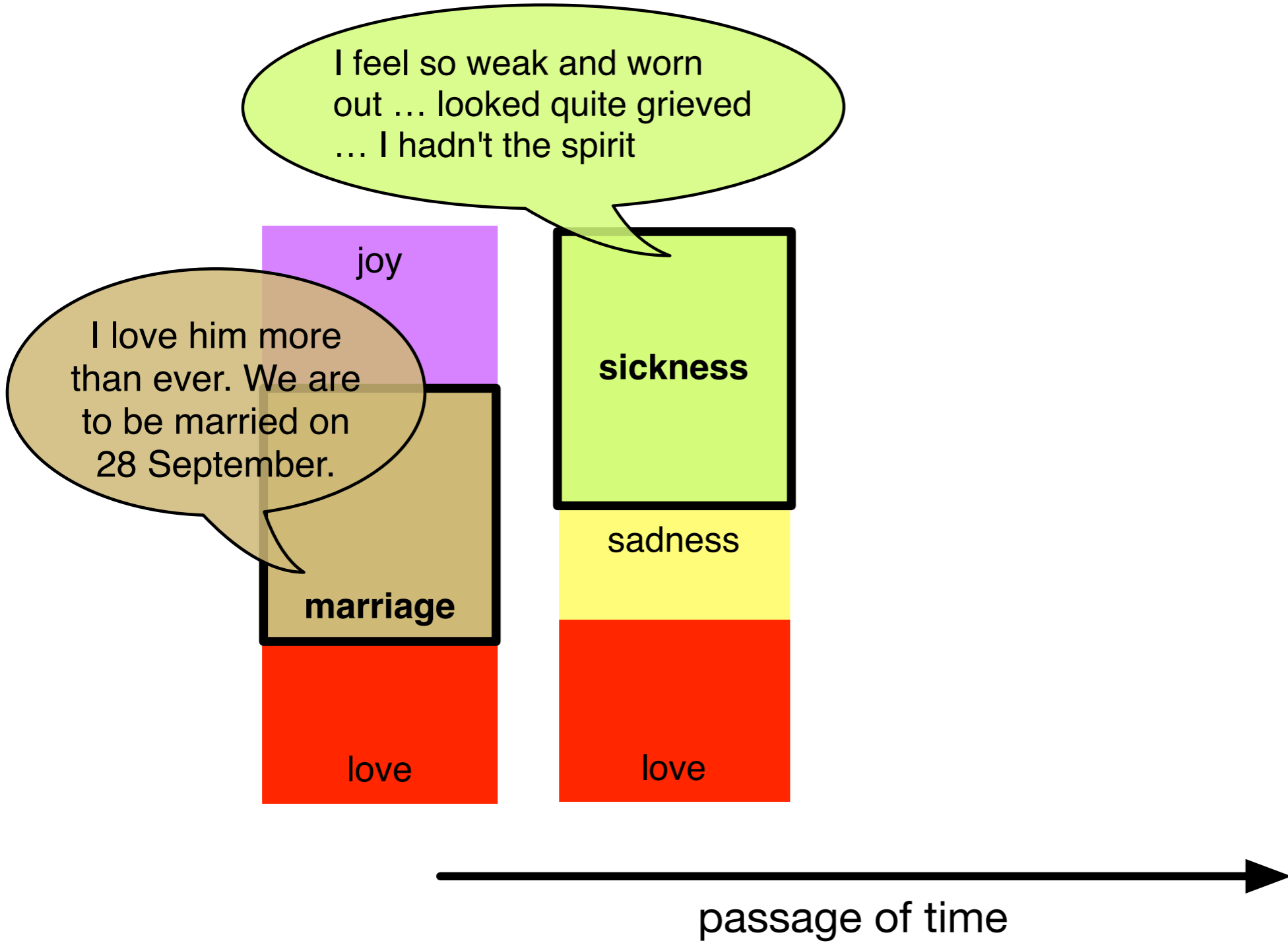
marriage

love

passage of time

Arthur and Lucy (*Dracula)*

joy

I feel so weak and worn out … looked quite grieved … I hadn't the spirit

I love him more than ever. We are to be married on 28 September.

marriage

love

passage of time

Arthur and Lucy (*Dracula)*

Arthur and Lucy (*Dracula*)

Arthur and Lucy (*Dracula*)

Arthur and Lucy (*Dracula*)

Arthur and Lucy (*Dracula*)

I love him more than ever. We are to be married on 28 September.

I feel so weak and worn out … looked quite grieved … I hadn't the spirit

poor girl, there is peace for her at last. It is the end!

Arthur placed the stake over her heart … he struck with all his might. The Thing in the coffin writhed …

joy

sickness

death

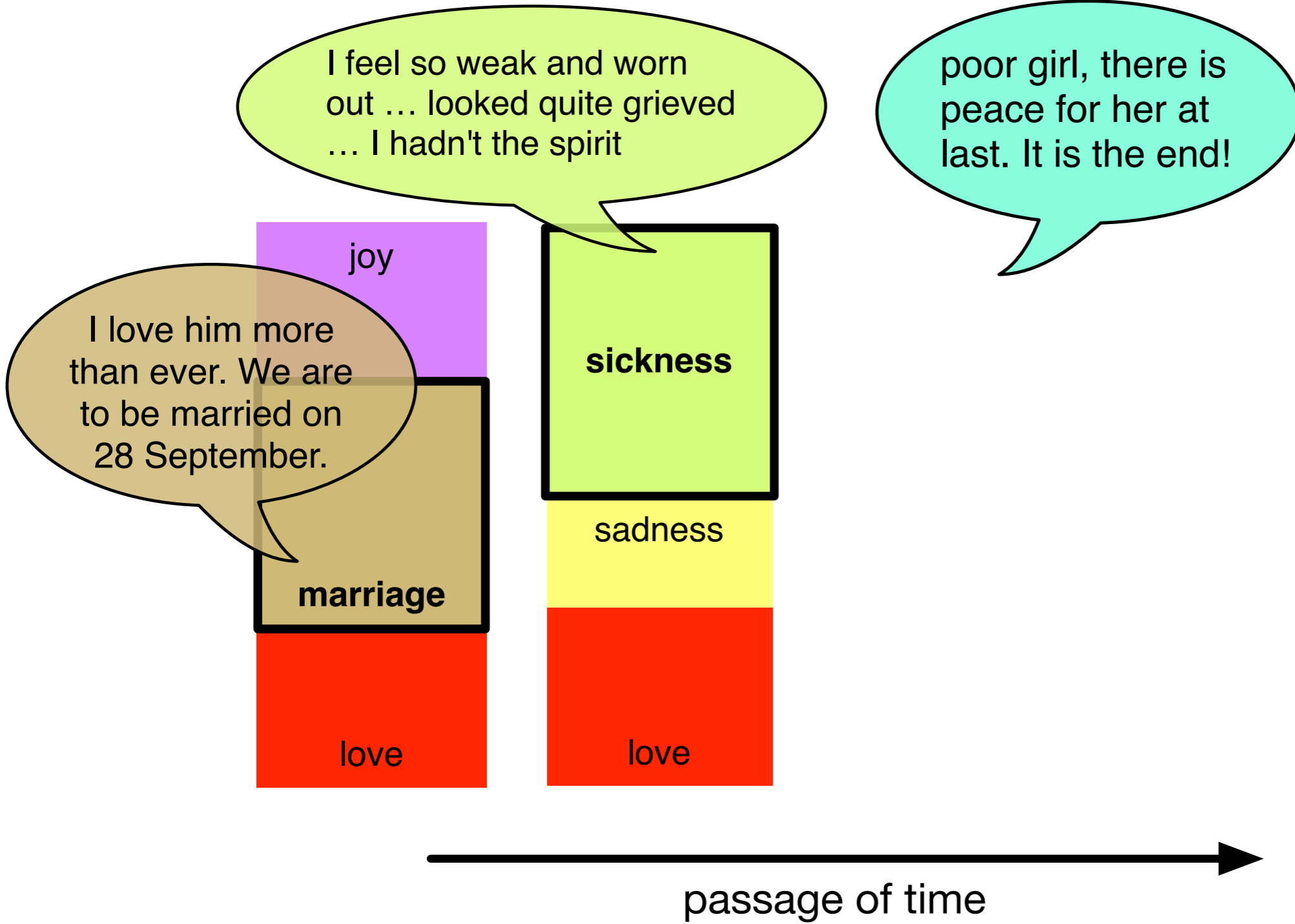death

fantasy

marriage

sadness

fantasy

sickness

murder

sadness

love

love

love

love

passage of time

Arthur and Lucy (*Dracula*)

# Why is this a worthwhile problem?

- "Distant reading" (Moretti, 2005) can help humanities scholars collect examples of specific relationship types

*"Do Jane Austen's female and male protagonists have a pattern in their evolving relationship (e.g., mutual disdain followed by romantic love)?"*
(Butler, 1975; Stovel, 1987; Hinant, 2006)

*"Do certain authors or novels portray relationships of desire more than others?"*
(Polhemus, 1990)

*"Can we detect positive or negative subtext underlying meals between two characters?"*
(Foster, 2009; Cognard-Black et al., 2014)

# Outline

- Dataset: character interactions

- RMN: relationship modeling network

- Experiments: coherent descriptors, interpretable trajectories

- Analysis: RMN's strengths and weaknesses

# A Dataset of Character Interactions

- For each pair of characters in a particular book, we extract all **spans** of text that contain mentions to both characters

# A Dataset of Character Interactions

- For each pair of characters in a particular book, we extract all **spans** of text that contain mentions to both characters

t=0

"If anyone was ever minding his business, it was I," Ignatius breathed. "Please. We must stop. I think I'm going to have a hemorrhage."
"Okay." Mrs. Reilly looked at her son's reddening face and realized that he would very happily collapse at her feet just to prove his point."

# A Dataset of Character Interactions

- For each pair of characters in a particular book, we extract all **spans** of text that contain mentions to both characters

t=0

"If anyone was ever minding his business, it was I," Ignatius breathed. "Please. We must stop. I think I'm going to have a hemorrhage." "Okay." Mrs. Reilly looked at her son's reddening face and realized that he would very happily collapse at her feet just to prove his point."

t=1

"Ignatius belched the gas of a dozen brownies trapped by his valve. "Grant me a little peace...." "You know I appreciate you, babe," Mrs. Reilly sniffed. "Come on and gimme a little goodbye kiss like a good boy."

# A Dataset of Character Interactions

- For each pair of characters in a particular book, we extract all **spans** of text that contain mentions to both characters

t=0

"If anyone was ever minding his business, it was I," Ignatius breathed. "Please. We must stop. I think I'm going to have a hemorrhage." "Okay." Mrs. Reilly looked at her son's reddening face and realized that he would very happily collapse at her feet just to prove his point."

t=1

"Ignatius belched the gas of a dozen brownies trapped by his valve. "Grant me a little peace…." "You know I appreciate you, babe," Mrs. Reilly sniffed. "Come on and gimme a little goodbye kiss like a good boy."

t=2

Mrs. Reilly looked at her son slyly and asked, "Ignatius, you sure you not a communiss?" "Oh, my God!" Ignatius bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

# A Dataset of Character Interactions

- 1,383 novels from Project Gutenberg and other Internet sources

  - Genres represented include romance, mystery, and fantasy

  - Preprocessed with David Bamman's BookNLP pipeline

  - Each span is a 200-token window centered around a character mention

- 20,013 unique character pairs and 380,408 spans

# Relationship Modeling Network (RMN)

- recurrent autoencoder with dictionary learning

# Relationship Modeling Network (RMN)

- recurrent autoencoder with dictionary learning



Mrs. Reilly looked at her son slyly and asked, "Ignatius, you sure you not a communiss?" "Oh, my God!" Ignatius bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

"Ignatius belched the gas of a dozen brownies trapped by his valve. "Grant me a little peace…." "You know I appreciate you, babe," Mrs. Reilly sniffed. "Come on and gimme a little goodbye kiss like a good boy."

Mrs. Reilly looked at her son slyly and asked, "Ignatius, you sure you not a communiss?" "Oh, my God!" Ignatius bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

# Relationship Modeling Network (RMN)

- recurrent autoencoder with dictionary learning

reconstruct inputs



Mrs. Reilly looked at her son slyly and asked, "Ignatius, you sure you not a communiss?" "Oh, my God!" Ignatius bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

"Ignatius belched the gas of a dozen brownies trapped by his valve. "Grant me a little peace…." "You know I appreciate you, babe," Mrs. Reilly sniffed. "Come on and gimme a little goodbye kiss like a good boy."

Mrs. Reilly looked at her son slyly and asked, "Ignatius, you sure you not a communiss?" "Oh, my God!" Ignatius bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

# Relationship Modeling Network (RMN)
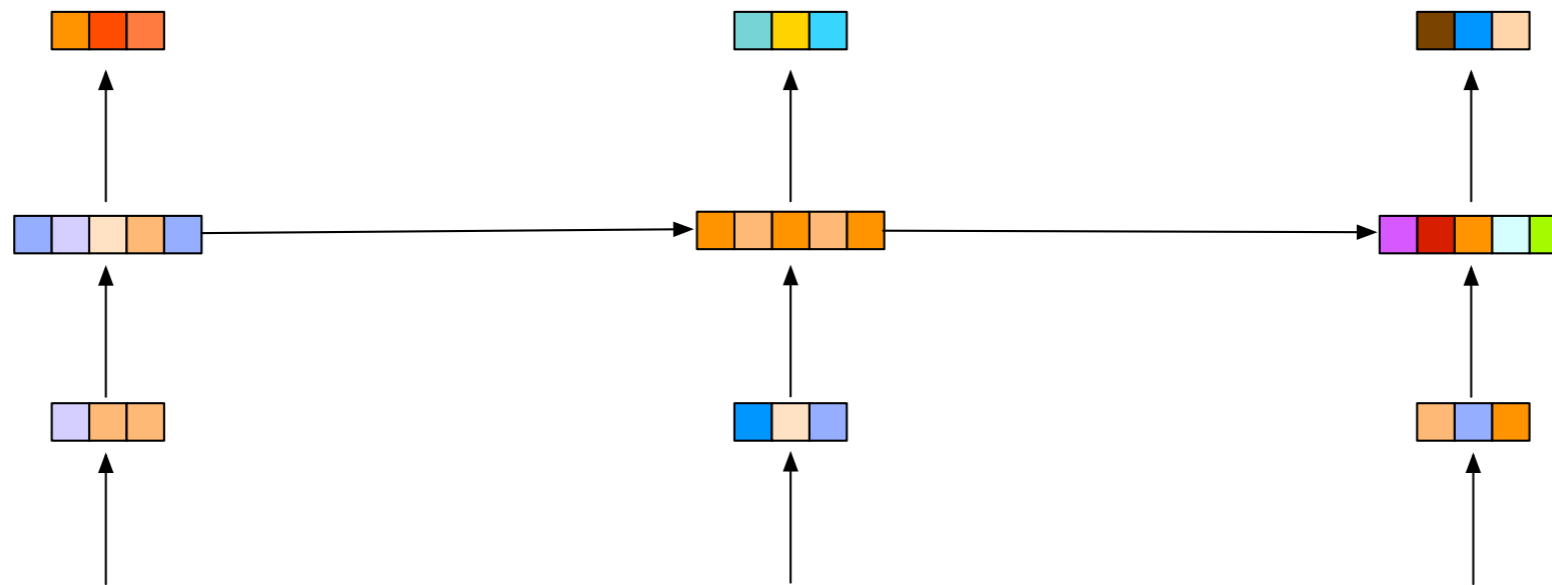
- recurrent autoencoder with dictionary learning



descriptor matrix

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

"**Ignatius** belched the gas of a dozen brownies trapped by his valve. "Grant me a little peace…." "You know I appreciate you, babe," **Mrs. Reilly** sniffed. "Come on and gimme a little goodbye kiss like a good boy."

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"
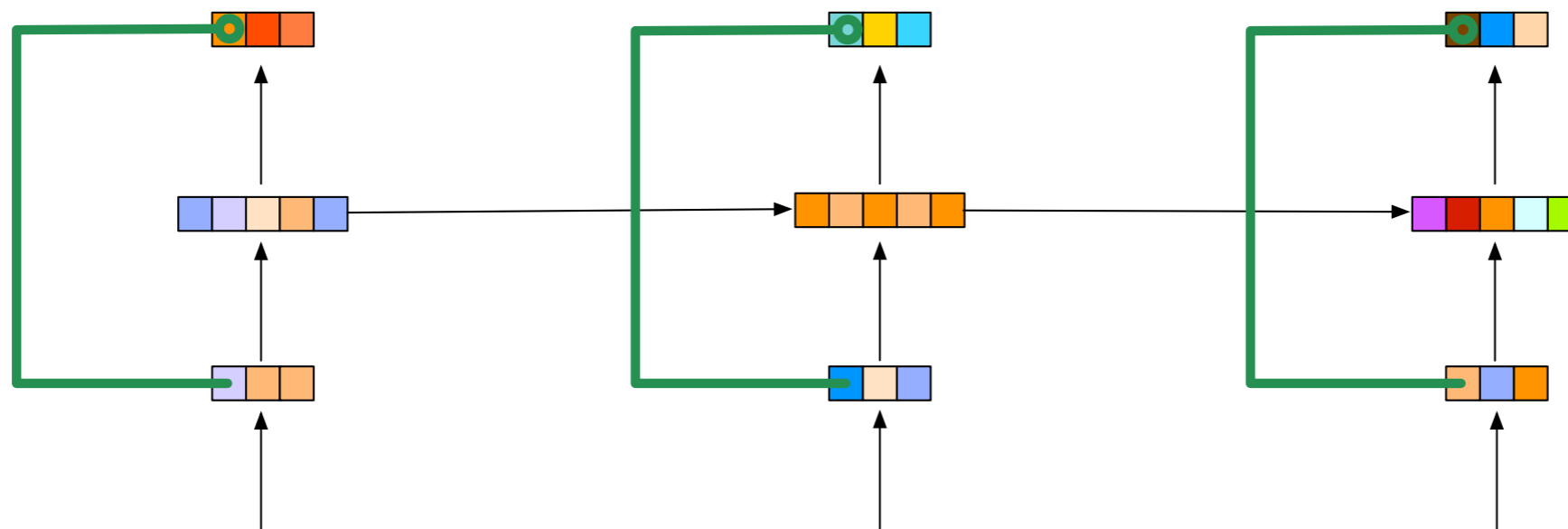
# Relationship Modeling Network (RMN)

- recurrent autoencoder with dictionary learning



reconstruct inputs

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

"**Ignatius** belched the gas of a dozen brownies trapped by his valve. "Grant me a little peace…." "You know I appreciate you, babe," **Mrs. Reilly** sniffed. "Come on and gimme a little goodbye kiss like a good boy."
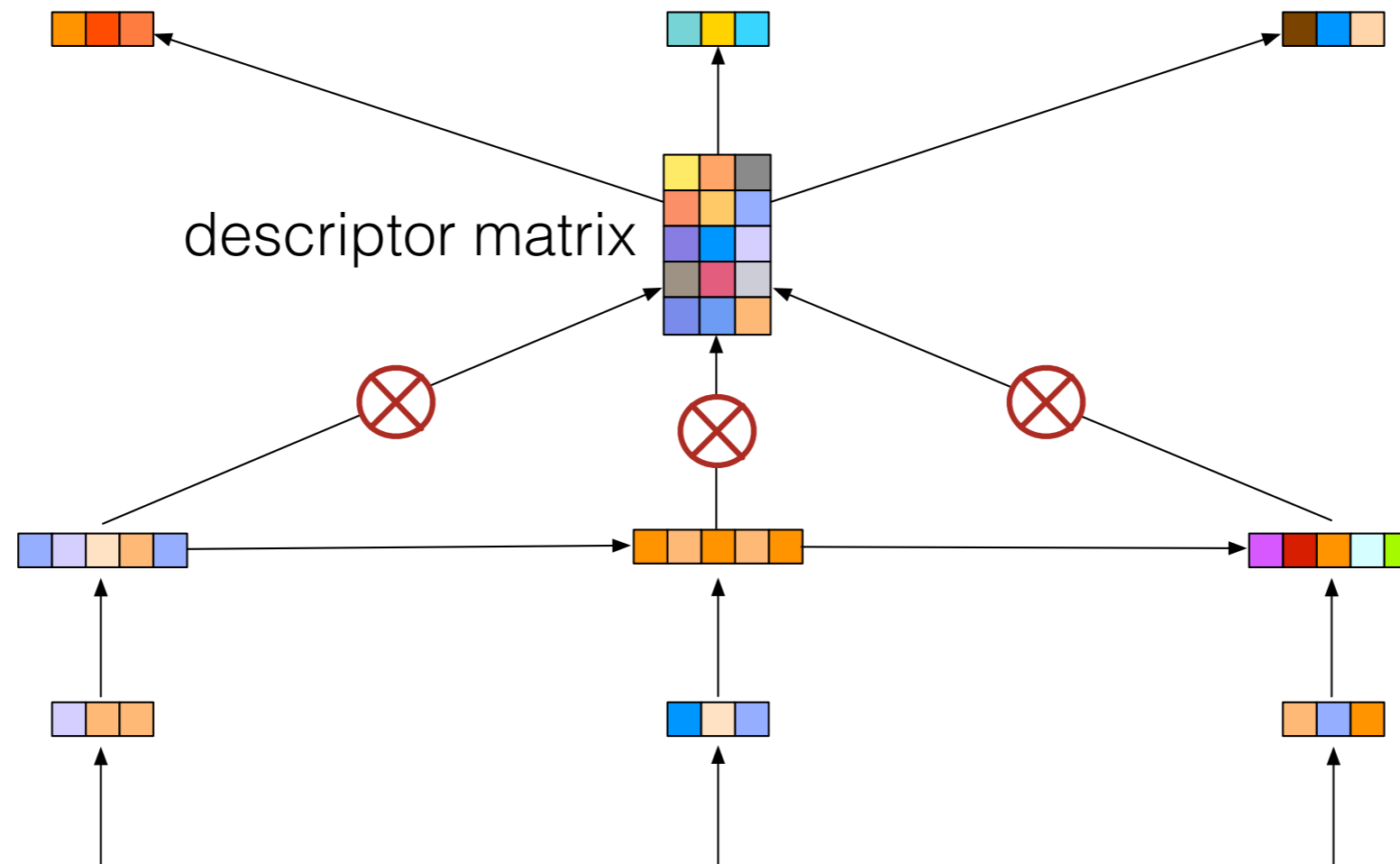
**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

$$\boldsymbol{v}_{s_t}$$

**1.** word embedding average

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

$$h_t = f(\mathbf{W}_h \cdot [\boldsymbol{v}_{s_t}; \boldsymbol{v}_{c_1}; \boldsymbol{v}_{c_2}; \boldsymbol{v}_b])$$

$\boldsymbol{v}_{s_t}$  $\boldsymbol{v}_{c_1}$  $\boldsymbol{v}_{c_2}$  $\boldsymbol{v}_b$

**2.** mix with embeddings for characters and books

**Mrs. Reilly**  **Ignatius**  "A Confederacy of Dunces"
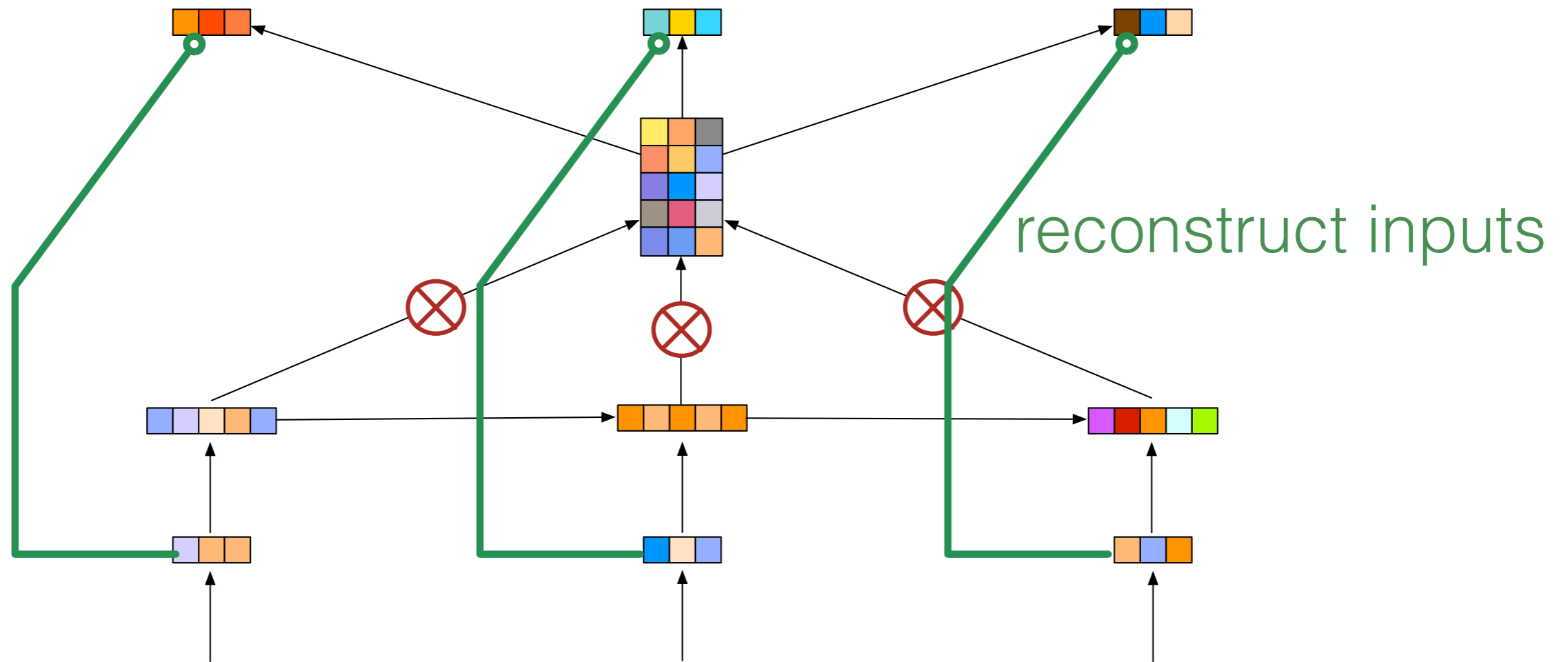
**1.** word embedding average

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

25

**4.** use a softmax activation function to make the hidden state a probability distribution over descriptors…

$$d_{t-1}: \textit{previous state}$$

**3.** a recurrent connection that copies over some of the previous hidden state

$$d_t = \alpha \cdot \mathrm{softmax}(\mathbf{W}_d \cdot [h_t; d_{t-1}]) +$$

$$(1 - \alpha) \cdot d_{t-1}: \textit{distribution over descriptors}$$

$$h_t = f(\mathbf{W}_h \cdot [v_{s_t}; v_{c_1}; v_{c_2}; v_b])$$

$$v_{s_t} \qquad v_{c_1} \qquad v_{c_2} \qquad v_b$$

**Mrs. Reilly**    **Ignatius**    "A Confederacy of Dunces"

**2.** mix with embeddings for characters and books
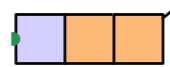
**1.** word embedding average

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

26

**4.** use a softmax activation function to make the hidden state a probability distribution over descriptors…

linear interpolation

$d_{t-1}$: *previous state*

**3.** a recurrent connection that copies over some of the previous hidden state

$$d_t = \alpha \cdot \mathrm{softmax}(\mathbf{W}_d \cdot [h_t; d_{t-1}]) +$$
$$(1 - \alpha) \cdot d_{t-1}: \text{\textit{distribution over descriptors}}$$

$$h_t = f(\mathbf{W}_h \cdot [v_{s_t}; v_{c_1}; v_{c_2}; v_b])$$

$v_{s_t}$　　$v_{c_1}$　　$v_{c_2}$　　$v_b$

**Mrs. Reilly**　　**Ignatius**　　"A Confederacy of Dunces"

**2.** mix with embeddings for characters and books

**1.** word embedding average

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

4. use a softmax
activation function
to make the hidden
state a probability
distribution over
descriptors…

softmax activation

$d_{t-1}$: *previous state*

**3.** a recurrent
connection that
copies over some
of the previous
hidden state

$$d_t = \alpha \cdot \text{softmax}(\mathbf{W}_d \cdot [h_t; d_{t-1}]) +$$

$$(1 - \alpha) \cdot d_{t-1}: \textit{distribution over}$$
*descriptors*

$$h_t = f(\mathbf{W}_h \cdot [v_{s_t}; v_{c_1}; v_{c_2}; v_b])$$

$v_{s_t}$     $v_{c_1}$     $v_{c_2}$     $v_b$

**Mrs. Reilly**    **Ignatius**    "A Confederacy
of Dunces"

**2.** mix with
embeddings for
characters and
books

**1.** word
embedding
average

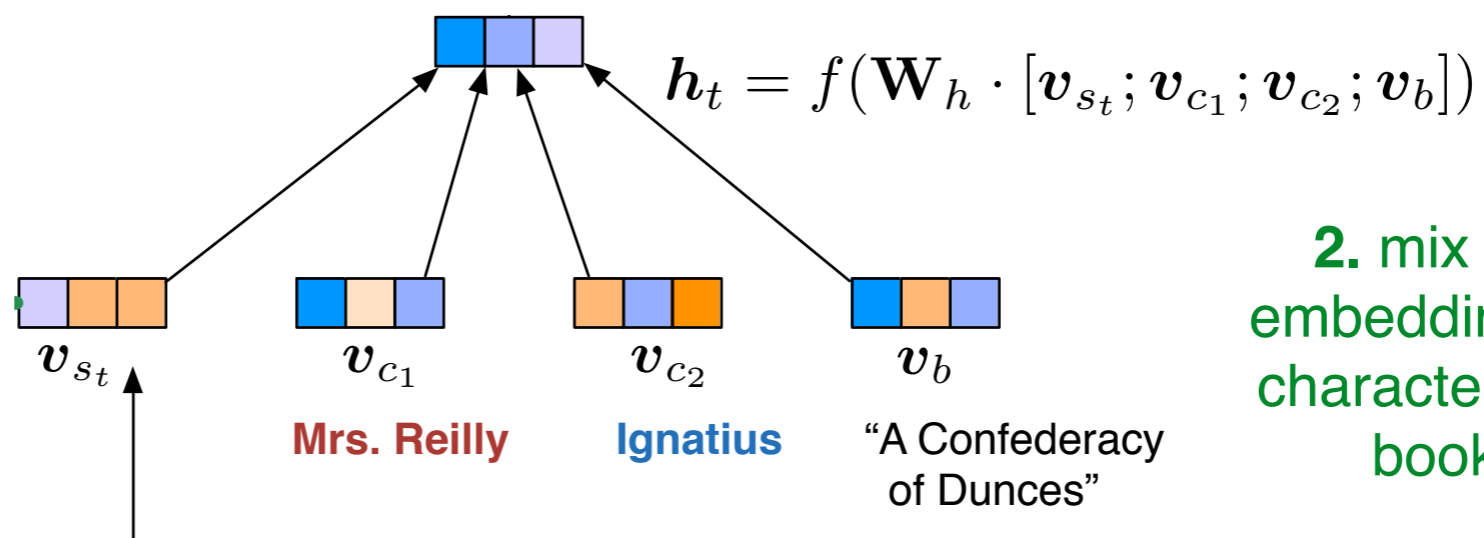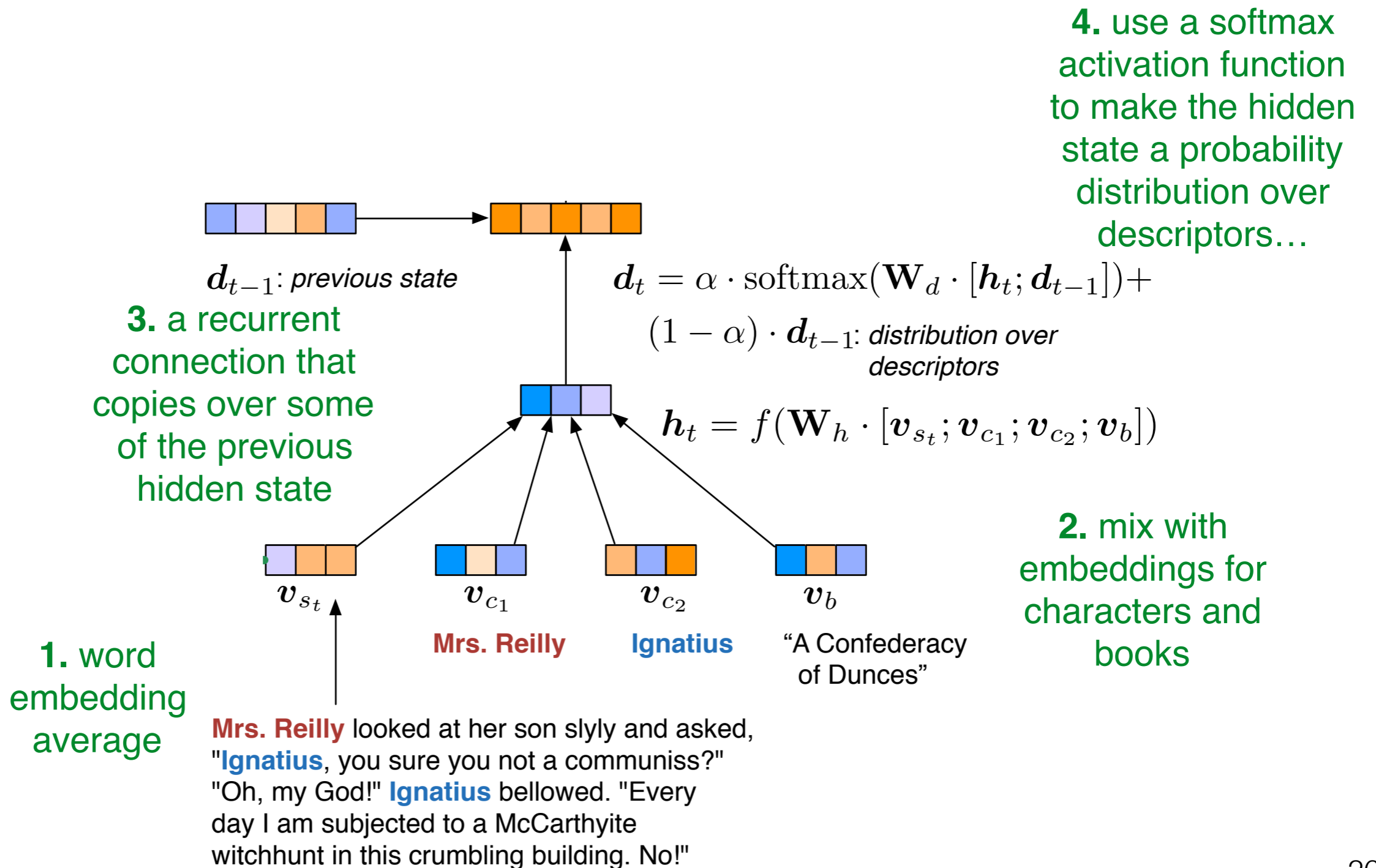**Mrs. Reilly** looked at her son slyly and asked,
"**Ignatius**, you sure you not a communiss?"
"Oh, my God!" **Ignatius** bellowed. "Every
day I am subjected to a McCarthyite
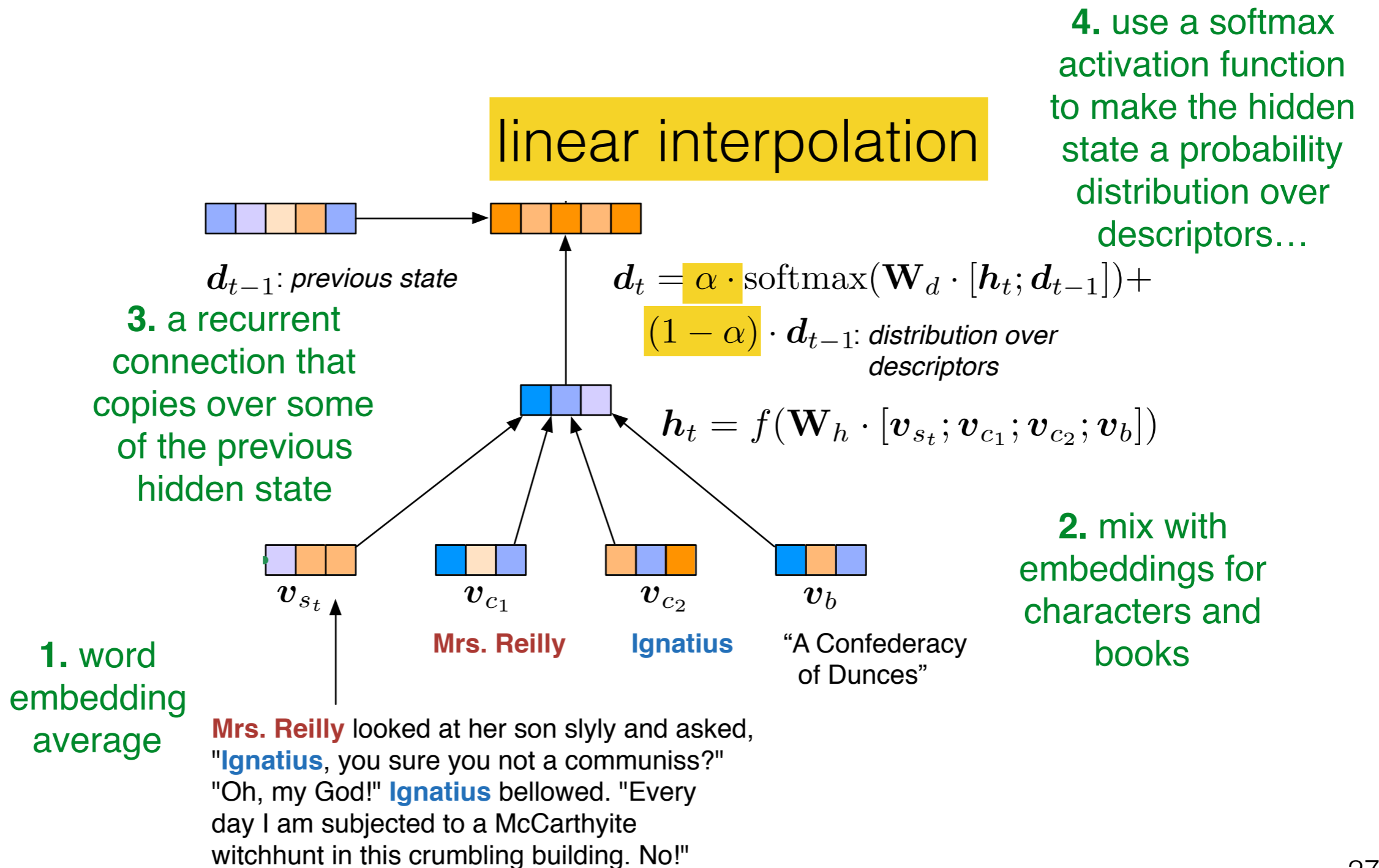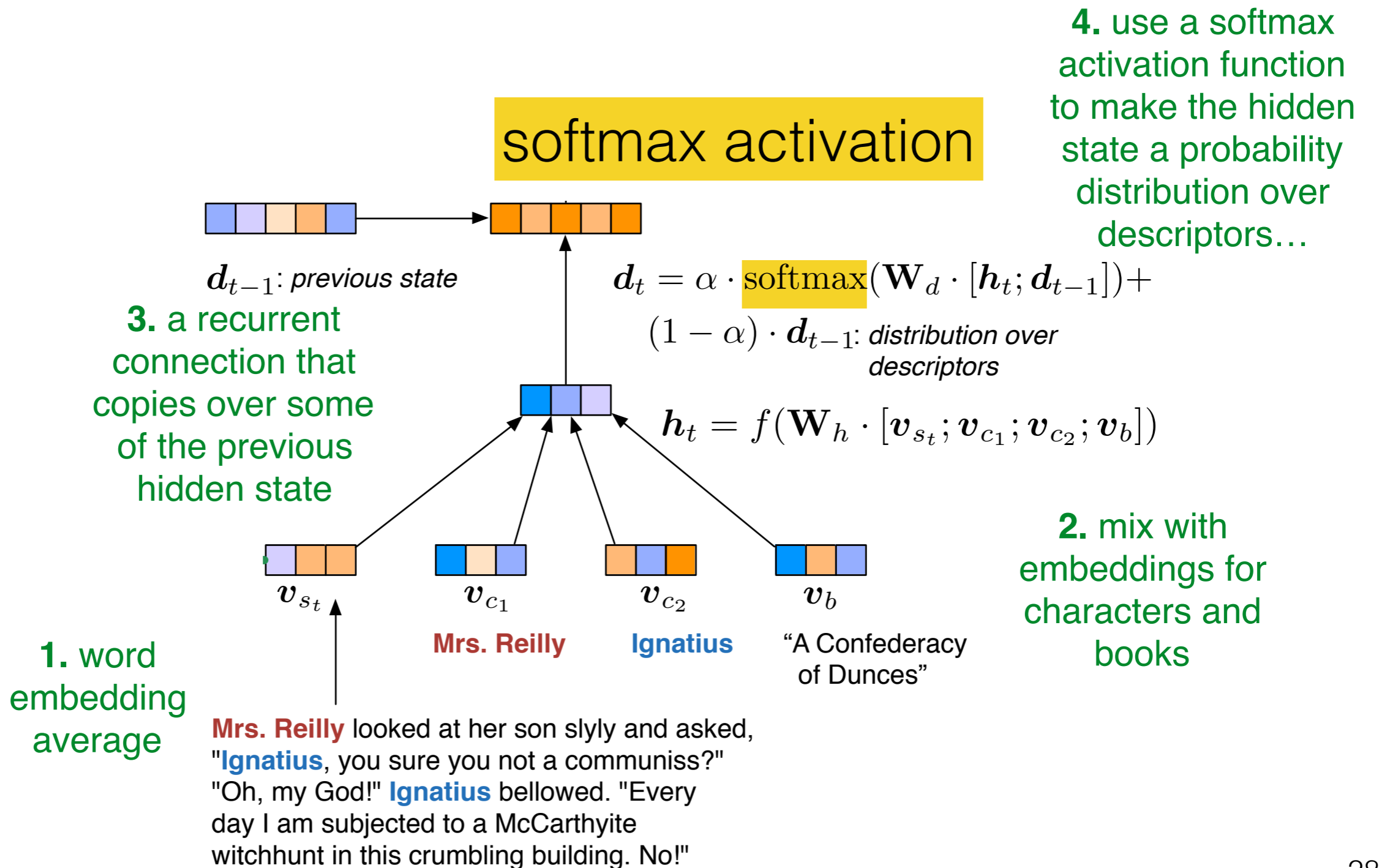witchhunt in this crumbling building. No!"

28

**5.** multiply the hidden state by the descriptor matrix to obtain a reconstruction of the span vector

$r_t = \mathbf{R}^\top d_t$ : *reconstruction of input span*

**4.** use a softmax activation function to make the hidden state a probability distribution over descriptors…

$\mathbf{R}$: *descriptor matrix*

$d_{t-1}$: *previous state*

$d_t = \alpha \cdot \mathrm{softmax}(\mathbf{W}_d \cdot [h_t; d_{t-1}]) +$

$(1 - \alpha) \cdot d_{t-1}$: *distribution over descriptors*

**3.** a recurrent connection that copies over some of the previous hidden state

$h_t = f(\mathbf{W}_h \cdot [v_{s_t}; v_{c_1}; v_{c_2}; v_b])$

**2.** mix with embeddings for characters and books

$v_{s_t}$       $v_{c_1}$       $v_{c_2}$       $v_b$

**Mrs. Reilly**       **Ignatius**       "A Confederacy of Dunces"

**1.** word embedding average

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

29

**6.** make the reconstructed vector close to the input span vector

**5.** multiply the hidden state by the descriptor matrix to obtain a reconstruction of the span vector

$r_t = \mathbf{R}^\top d_t$ : *reconstruction of input span*

**4.** use a softmax activation function to make the hidden state a probability distribution over descriptors…

$\mathbf{R}$: *descriptor matrix*

$d_{t-1}$: *previous state*

**3.** a recurrent connection that copies over some of the previous hidden state

$d_t = \alpha \cdot \text{softmax}(\mathbf{W}_d \cdot [h_t; d_{t-1}]) +$
$(1 - \alpha) \cdot d_{t-1}$: *distribution over descriptors*
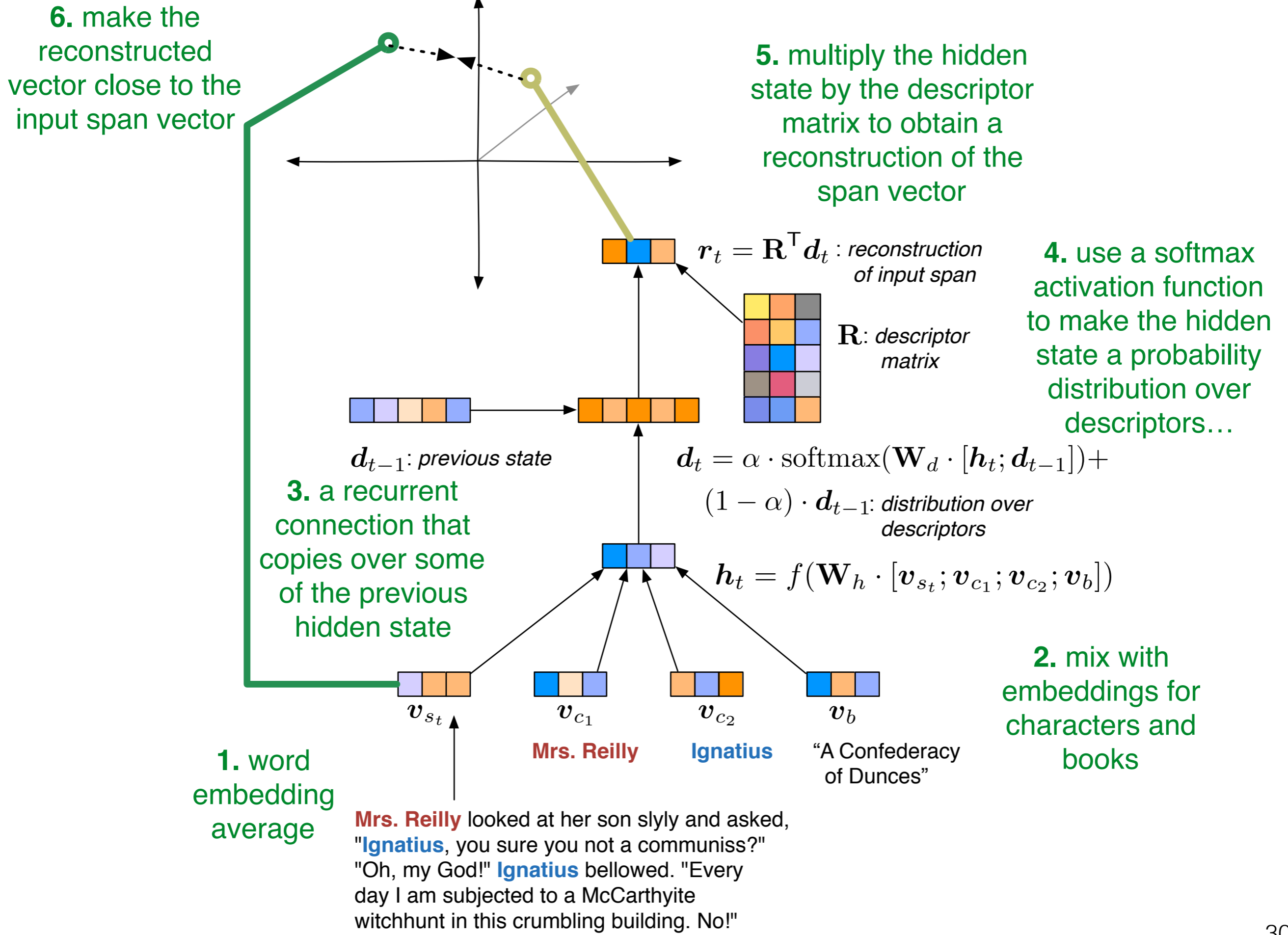
$h_t = f(\mathbf{W}_h \cdot [v_{s_t}; v_{c_1}; v_{c_2}; v_b])$

**2.** mix with embeddings for characters and books

$v_{s_t}$

$v_{c_1}$

$v_{c_2}$

$v_b$

**Mrs. Reilly**

**Ignatius**

"A Confederacy of Dunces"

**1.** word embedding average

**Mrs. Reilly** looked at her son slyly and asked, "**Ignatius**, you sure you not a communiss?" "Oh, my God!" **Ignatius** bellowed. "Every day I am subjected to a McCarthyite witchhunt in this crumbling building. No!"

30

# Labeling the Learned Descriptors

- We compute the nearest word embeddings to each row of the descriptor matrix **R**, which humans use to provide external labels.
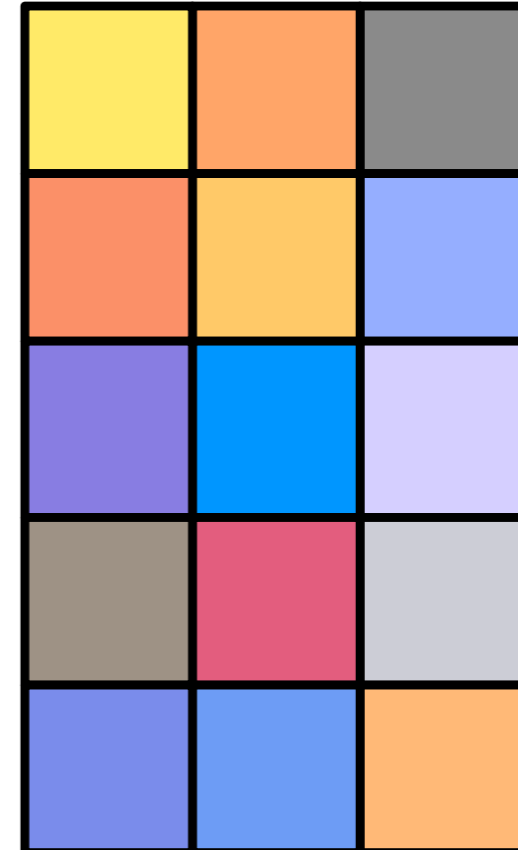
**violence:** *grenades, guns, bullets*

**sadness:** *regretful, rueful, pity*

**politics:** *political, leadership, rule*

**fantasy:** *cosmic, realms, universe*

**suffering:** *fear, nightmares, suffer*

# Relationship to Topic Models

- RMN outputs ≈ topic model latent variables:

  - descriptor matrix **R** ≈ topic-word matrices **φ**

  - descriptor weights $d_t$ at each timestep ≈ document-topic assignments **z**

- Baselines:

  - temporally-oblivious: **LDA** (Blei et al., 2001), **Nubbi** (Chang et al., 2008)

  - temporally-aware: **HTMM** (Gruber et al., 2007)

# Experiment 1:
Descriptor Coherence

# Do the Descriptors Make Sense?

- Goal: compare the descriptors learned by the RMN to the topics learned by our topic model baselines

- Task: word intrusion (Chang et al., 2009)

  - Workers identify an "intruder" word from a set of words that —other than the intruder—come from the same descriptor

*contempt malice condescend praise distaste mock*

*worship pray devote yourselves gods gather*

# Do the Descriptors Make Sense?

- Goal: compare the descriptors learned by the RMN to the topics learned by our topic model baselines

- Task: word intrusion (Chang et al., 2009)

    - Workers identify an "intruder" word from a set of words that —other than the intruder—come from the same descriptor

*contempt malice condescend **praise** distaste mock*

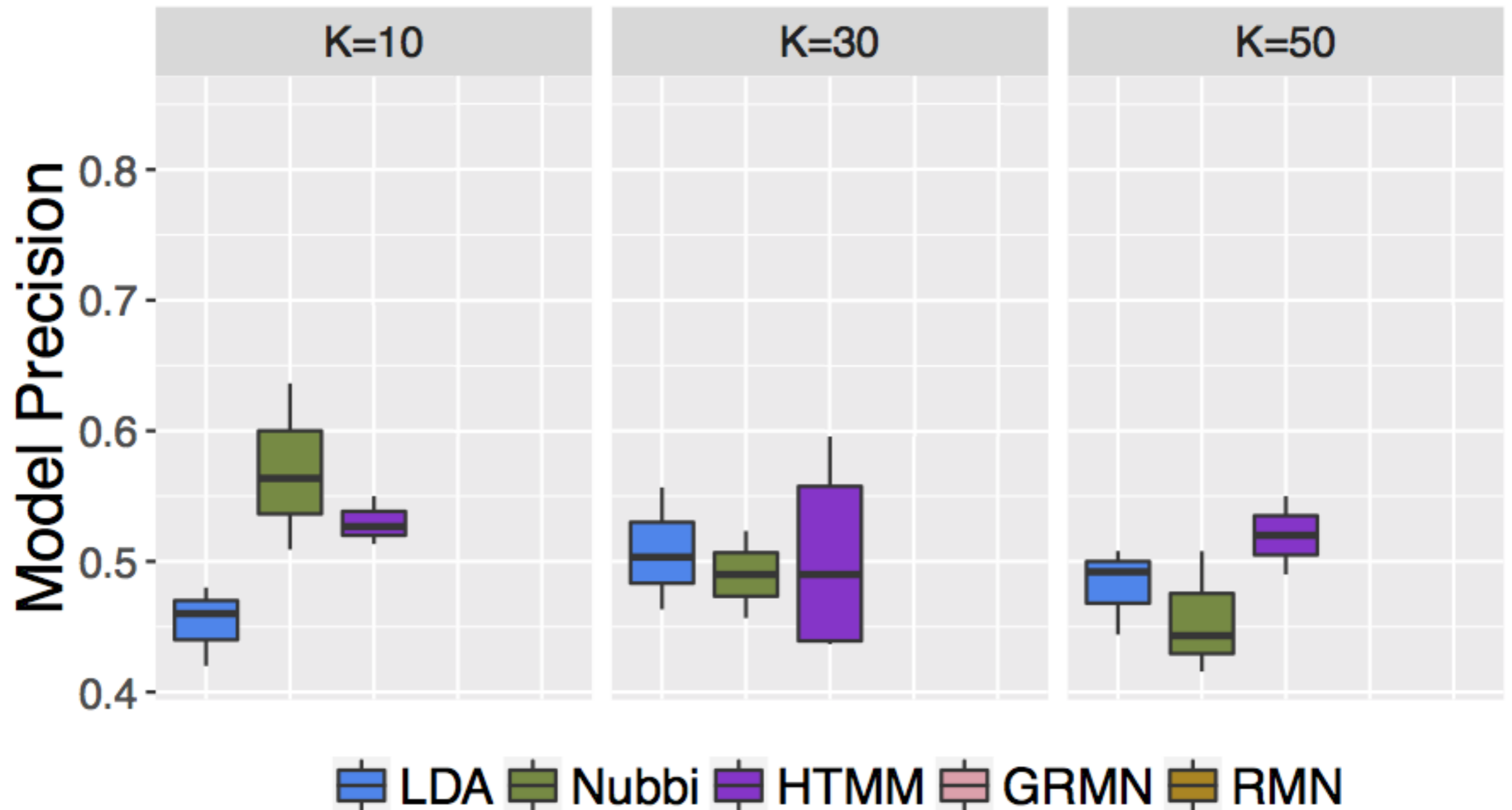*worship pray devote yourselves gods gather*

# Do the Descriptors Make Sense?

- Goal: compare the descriptors learned by the RMN to the topics learned by our topic model baselines

- Task: word intrusion (Chang et al., 2009)

  - Workers identify an "intruder" word from a set of words that —other than the intruder—come from the same descriptor
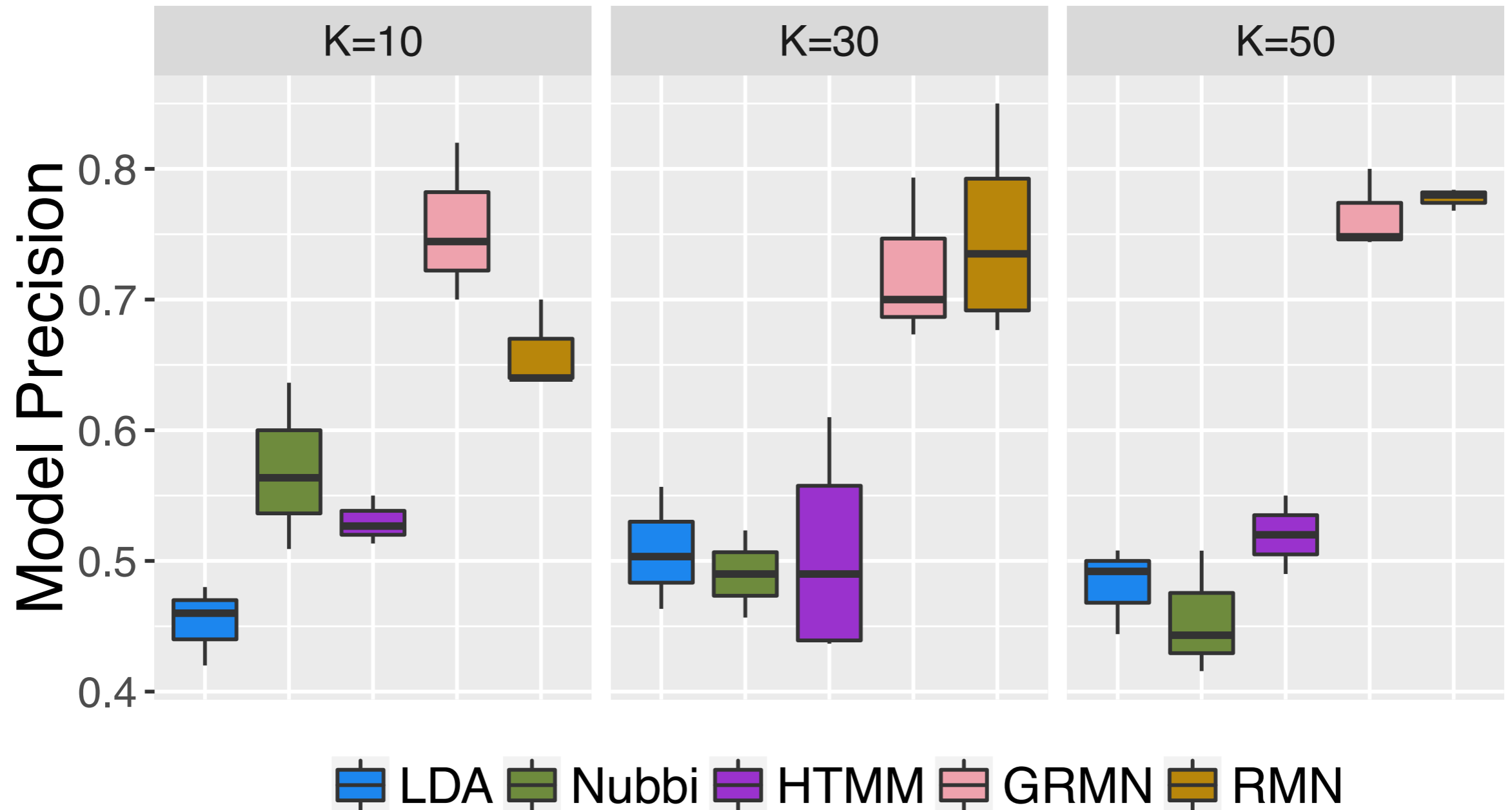
*contempt malice condescend **praise** distaste mock*

*worship pray devote **yourselves** gods **gather***

# Do the Descriptors Make Sense?

# Do the Descriptors Make Sense?

# Coherent Descriptors

## RMN

**outdoors:** *outdoors trail trails hillside grassy slopes*
**sadness:** *regretful rueful pity pained despondent*
**education:** *teaching graduate year teacher attended*
**love:** *love delightful happiness enjoyed enjoyable*
**murder:** *autopsy arrested homicide murdered*

## HTMM

**crime:** *blood knife pain legs steal*
**food:** *kitchen mouth glass food bread*
**violence:** *sword shot blood shouted swung*
**boats:** *ship boat captain deck crew*
**outdoors:** *stone rock path darkness desert*

# Coherent Descriptors

**RMN**

**outdoors:** *outdoors trail trails hillside grassy slopes*
**sadness:** *regretful rueful pity pained despondent*
**education:** *teaching graduate year teacher attended*
**love:** *love delightful happiness enjoyed enjoyable*
**murder:** *autopsy arrested homicide murdered*

**HTMM**

**crime:** *blood knife pain legs steal*
**food:** *kitchen mouth glass food bread*
**violence:** *sword shot blood shouted swung*
**boats:** *ship boat captain deck crew*
**outdoors:** *stone rock path darkness desert*

# Experiment 2:
## Trajectory Quality

# Visualizing Trajectories

- for all time steps $t$, compute *argmax* of $d_t$ and stack vertically

# Visualizing Trajectories

- for all time steps *t*, compute *argmax* of $d_t$ and stack vertically

| time | love | death | money | crime |
|------|------|-------|-------|-------|
| 0 | 0.95 | 0.01 | 0.03 | 0.01 |
| 1 | 0.8 | 0.01 | 0.18 | 0.01 |
| 2 | 0.4 | 0.01 | 0.5 | 0.09 |
| 3 | 0.3 | 0.01 | 0.2 | 0.5 |
| 4 | 0.2 | 0.7 | 0.05 | 0.05 |

# Visualizing Trajectories

- for all time steps $t$, compute *argmax* of $d_t$ and stack vertically

| time | love | death | money | crime |
|------|------|-------|-------|-------|
| 0 | **0.95** | 0.01 | 0.03 | 0.01 |
| 1 | 0.8 | 0.01 | 0.18 | 0.01 |
| 2 | 0.4 | 0.01 | 0.5 | 0.09 |
| 3 | 0.3 | 0.01 | 0.2 | 0.5 |
| 4 | 0.2 | 0.7 | 0.05 | 0.05 |

**time**

love

# Visualizing Trajectories

- for all time steps $t$, compute *argmax* of $d_t$ and stack vertically

| time | love | death | money | crime |
|------|------|-------|-------|-------|
| 0 | **0.95** | 0.01 | 0.03 | 0.01 |
| 1 | **0.8** | 0.01 | 0.18 | 0.01 |
| 2 | 0.4 | 0.01 | 0.5 | 0.09 |
| 3 | 0.3 | 0.01 | 0.2 | 0.5 |
| 4 | 0.2 | 0.7 | 0.05 | 0.05 |

**time**

love

# Visualizing Trajectories

- for all time steps $t$, compute *argmax* of $d_t$ and stack vertically

| time | love | death | money | crime |
|------|------|-------|-------|-------|
| 0 | **0.95** | 0.01 | 0.03 | 0.01 |
| 1 | **0.8** | 0.01 | 0.18 | 0.01 |
| 2 | 0.4 | 0.01 | **0.5** | 0.09 |
| 3 | 0.3 | 0.01 | 0.2 | 0.5 |
| 4 | 0.2 | 0.7 | 0.05 | 0.05 |



time

money

love

# Visualizing Trajectories

- for all time steps $t$, compute *argmax* of $d_t$ and stack vertically

| time | love | death | money | crime |
|------|------|-------|-------|-------|
| 0 | **0.95** | 0.01 | 0.03 | 0.01 |
| 1 | **0.8** | 0.01 | 0.18 | 0.01 |
| 2 | 0.4 | 0.01 | **0.5** | 0.09 |
| 3 | 0.3 | 0.01 | 0.2 | **0.5** |
| 4 | 0.2 | 0.7 | 0.05 | 0.05 |



crime
money
love

# Visualizing Trajectories

- for all time steps *t*, compute *argmax* of $d_t$ and stack vertically

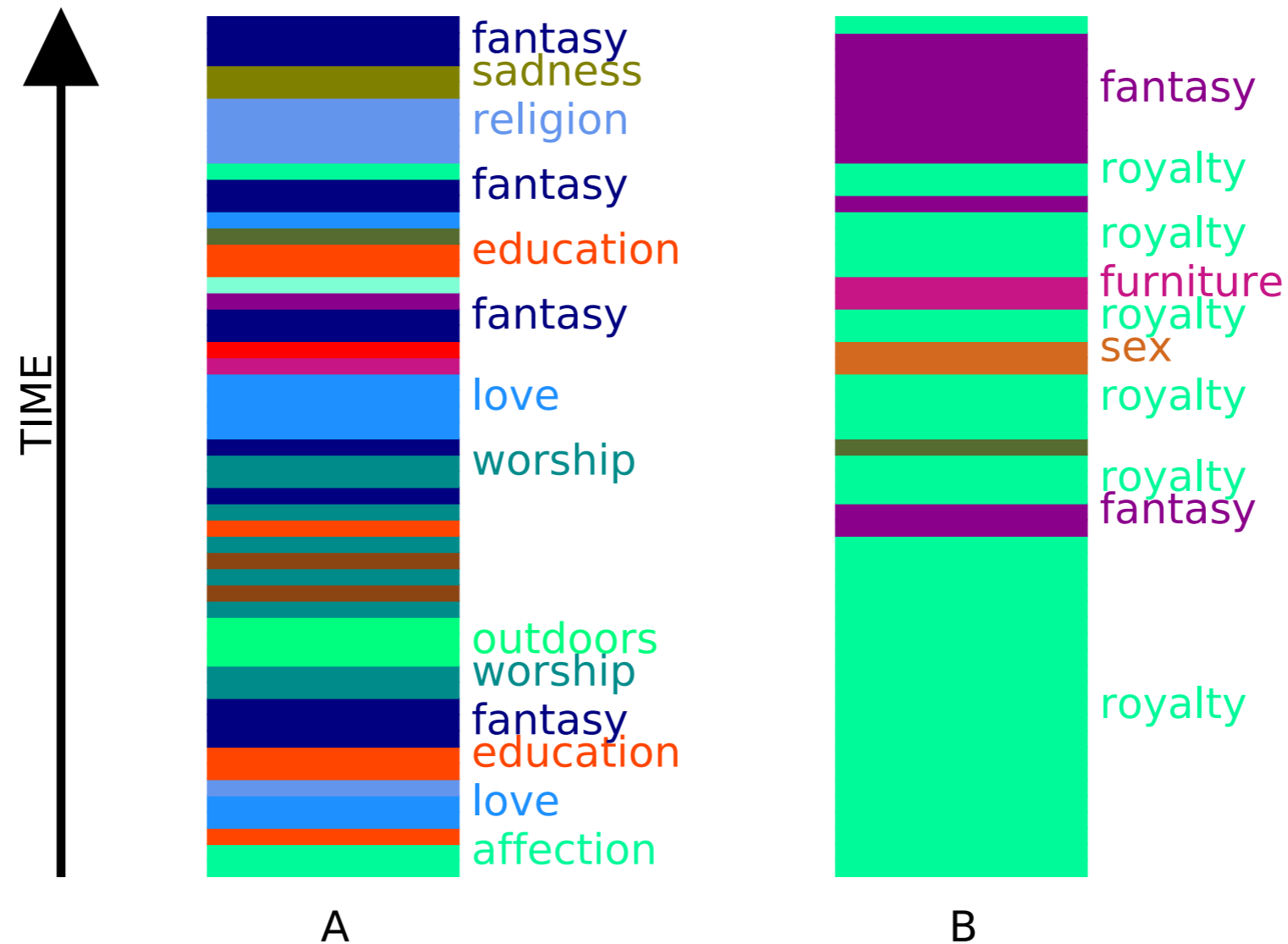| time | love | death | money | crime |
|------|------|-------|-------|-------|
| 0 | **0.95** | 0.01 | 0.03 | 0.01 |
| 1 | **0.8** | 0.01 | 0.18 | 0.01 |
| 2 | 0.4 | 0.01 | **0.5** | 0.09 |
| 3 | 0.3 | 0.01 | 0.2 | **0.5** |
| 4 | 0.2 | **0.7** | 0.05 | 0.05 |



death
crime
money
love

time

# Do the Trajectories Make Sense?

In this task, you will be comparing two timelines of how a relationship between a pair of literary characters changes over time. We will provide you with a summary of the relationship, and your job is to select which of the two timelines (A or B) better captures the content of the summary.

- We crawl Wikipedia and SparkNotes for summaries

- Removing uninformative summaries results in 125 character pairs to evaluate

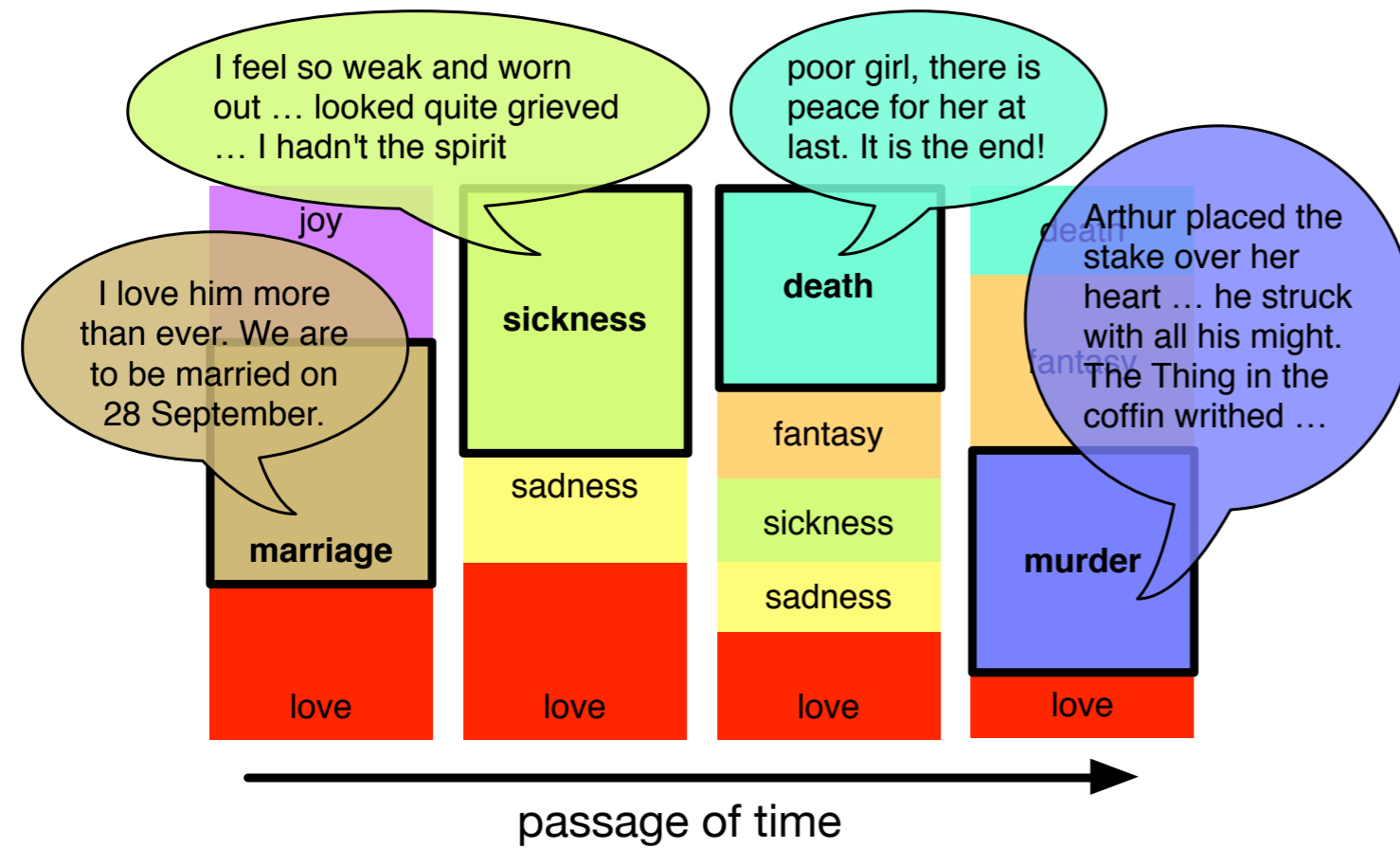- Workers prefer the **RMN** to the **HTMM** for 87 out of the 125 relationships (69.6%, Fleiss κ=0.32)

Siddhartha: Siddhartha AND Govinda

A

B

TIME

Column A (top to bottom): fantasy, sadness, religion, fantasy, education, fantasy, love, worship, outdoors, worship, fantasy, education, love, affection

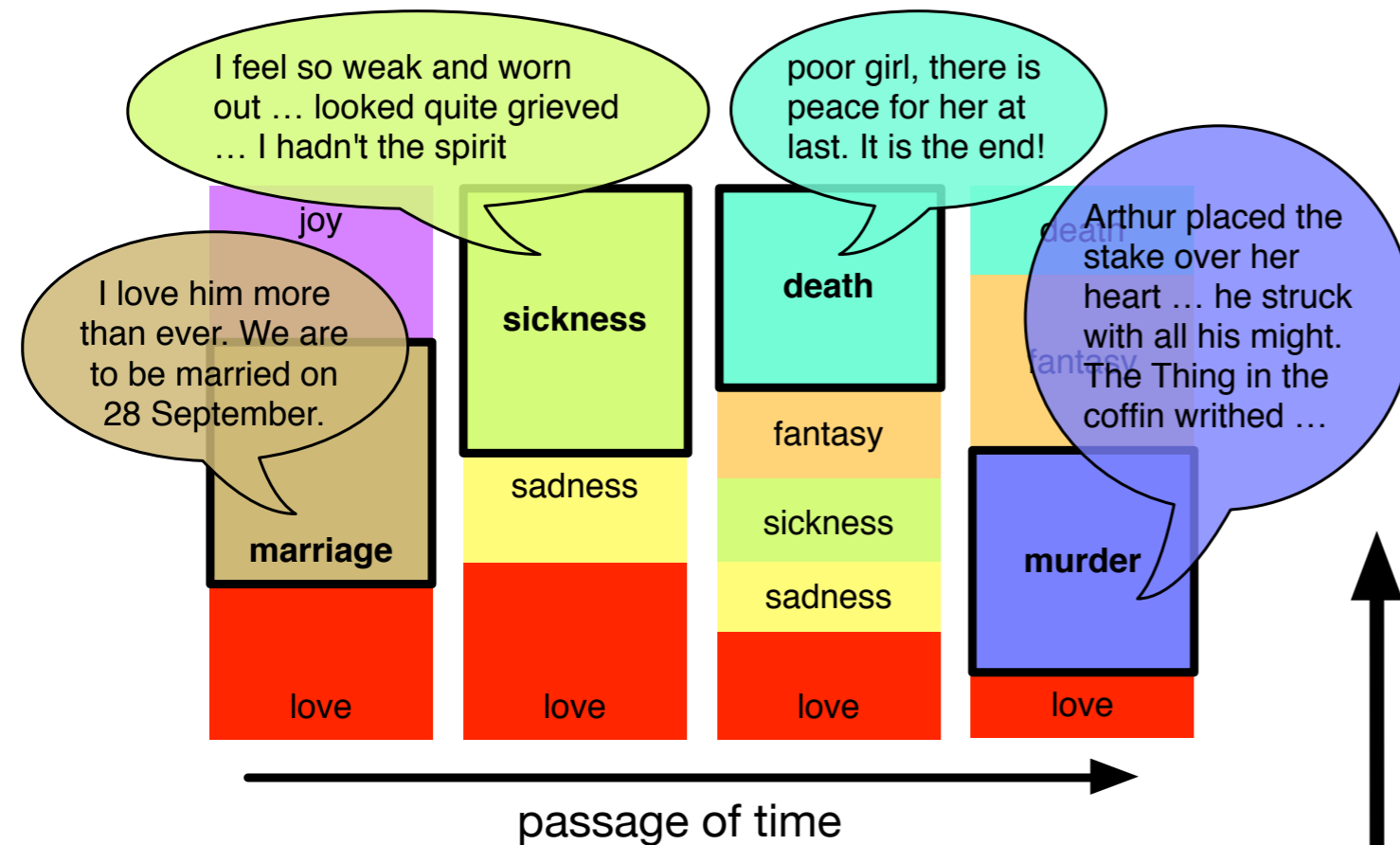Column B (top to bottom): fantasy, royalty, royalty, furniture, royalty, sex, royalty, royalty, fantasy, royalty

**Summary: Govinda** is **Siddhartha**'s best friend and sometimes his follower. Like **Siddhartha**, **Govinda** devotes his life to the quest for understanding and enlightenment. He leaves his village with **Siddhartha** to join the Samanas, then leaves the Samanas to follow Gotama. He searches for enlightenment independently of **Siddhartha** but persists in looking for teachers who can show him the way. In the end, he is able to achieve enlightenment only because of **Siddhartha**'s love for him.

# Qualitative Analysis:
Good and Bad Trajectories

# Arthur and Lucy "ground-truth":
## marriage -> sickness -> death -> murder



passage of time

Arthur and Lucy "ground-truth":
marriage -> sickness -> death -> murder

learned trajectories:

**Storm Island**: **David** and **Lucy**

RMN        HTMM

Event-based similarities
between the two models

**Storm Island**: David and Lucy — RMN / HTMM

**A Tale of Two Cities**: Darnay and Lucie — RMN / HTMM

Event-based similarities between the two models

The RMN is led astray by the novel's sad tone

55

# Qualitative Analysis:
Using Existing Datasets

# What Makes a Good Relationship?

- Dataset of Massey et al. (2015) has affinity annotations for relationships in Project Gutenberg

  - 120 non-neutral relationships are also present in our dataset

# What Makes a Good Relationship?

- Dataset of Massey et al. (2015) has affinity annotations for relationships in Project Gutenberg

  - 120 non-neutral relationships are also present in our dataset

<div>

**positive**

| love | death | sadness |
|------|-------|---------|
| 0.9 | 0.05 | 0.05 |
| 0.8 | 0.1 | 0.1 |
| 0.6 | 0.3 | 0.1 |
| 0.7 | 0.1 | 0.2 |
| 0.8 | 0.1 | 0.1 |

Don Quixote & Sancho Panza
Candide & Cunégonde
Anna Karenina & Vronsky
…

</div>

<div>

**negative**

| love | death | sadness |
|------|-------|---------|
| 0.1 | 0.7 | 0.2 |
| 0.2 | 0.3 | 0.5 |
| 0.15 | 0.25 | 0.6 |
| 0.05 | 0.65 | 0.3 |
| 0.1 | 0.2 | 0.7 |

Dracula & Jonathan Harker
Dr. Jekyll & Mr. Hyde
Hester Prynne & Chillingworth
…

</div>

# What Makes a Good Relationship?

- Dataset of Massey et al. (2015) has affinity annotations for relationships in Project Gutenberg

  - 120 non-neutral relationships are also present in our dataset

positive

| love | death | sadness |
|------|-------|---------|
| 0.9  | 0.05  | 0.05    |
| 0.8  | 0.1   | 0.1     |
| 0.6  | 0.3   | 0.1     |
| 0.7  | 0.1   | 0.2     |
| 0.8  | 0.1   | 0.1     |
| **0.76** | **0.13** | **0.11** |

negative

| love | death | sadness |
|------|-------|---------|
| 0.1  | 0.7   | 0.2     |
| 0.2  | 0.3   | 0.5     |
| 0.15 | 0.25  | 0.6     |
| 0.05 | 0.65  | 0.3     |
| 0.1  | 0.2   | 0.7     |
| **0.12** | **0.42** | **0.46** |

average the positive and negative trajectories

# What Makes a Good Relationship?

- Dataset of Massey et al. (2015) has affinity annotations for relationships in Project Gutenberg
  - 120 non-neutral relationships are also present in our dataset

<table>
<tr><th colspan="3">positive</th></tr>
<tr><th>love</th><th>death</th><th>sadness</th></tr>
<tr><td>0.9</td><td>0.05</td><td>0.05</td></tr>
<tr><td>0.8</td><td>0.1</td><td>0.1</td></tr>
<tr><td>0.6</td><td>0.3</td><td>0.1</td></tr>
<tr><td>0.7</td><td>0.1</td><td>0.2</td></tr>
<tr><td>0.8</td><td>0.1</td><td>0.1</td></tr>
<tr><td><b>0.76</b></td><td><b>0.13</b></td><td><b>0.11</b></td></tr>
</table>

1. love
2. death
3. sadness

<table>
<tr><th colspan="3">negative</th></tr>
<tr><th>love</th><th>death</th><th>sadness</th></tr>
<tr><td>0.1</td><td>0.7</td><td>0.2</td></tr>
<tr><td>0.2</td><td>0.3</td><td>0.5</td></tr>
<tr><td>0.15</td><td>0.25</td><td>0.6</td></tr>
<tr><td>0.05</td><td>0.65</td><td>0.3</td></tr>
<tr><td>0.1</td><td>0.2</td><td>0.7</td></tr>
<tr><td><b>0.12</b></td><td><b>0.42</b></td><td><b>0.46</b></td></tr>
</table>

1. sadness
2. death
3. love

# Most Positive Descriptors

**RMN**

education
love
religion
sex

**HTMM**

love
parental
business
outdoors

# Most Positive Descriptors

**RMN**
education
love
religion
sex

**HTMM**
love
parental
business
outdoors

# Most Negative Descriptors

**RMN**
politics
murder
sadness
royalty

**HTMM**
love
politics
violence
crime

Most Positive Descriptors

RMN
education
love
religion
sex

HTMM
love
parental
business
outdoors

Most Negative Descriptors

RMN
**politics**
murder
sadness
royalty

HTMM
love
**politics**
violence
crime

# Why is Politics Negative?

- Both models rank **politics** as highly negative

- The affinity data we look at comes primarily from Victorian-era authors (e.g., Charles Dickens and George Eliot)

*Victorian-era authors are "obsessed with otherness… of antiquated social and legal institutions, and of autocratic and/or dictatorial abusive government"* (Zarifopol-Johnston,1995)

# Areas for Improvement

- Difficult to evaluate unsupervised relationship modeling, requires considerable human effort

- Our data processing leaves out a lot of information

  - e.g., spans of text in which one but not both characters in a relationship are mentioned

  - only considers *undirected* relationships between *pairs*

- Model performance is directly tied to the quality of character disambiguation and coreference resolution

  - e.g., first person pronouns

# Recap

- Introduced the task of unsupervised relationship modeling as well as an *interpretable* neural network architecture, the RMN, for this task

- Found that the RMN generates higher quality descriptors and more interpretable trajectories than topic model baselines

- Future work: collaborate with humanities researchers to help answer literary questions with the RMN

# Thanks! Questions?

code/data @ github.com/miyyer/rmn