# Generative Models of Monolingual and Bilingual Gappy Patterns

Kevin Gimpel          Noah A. Smith

# Overview

- We present models that generate text using patterns with gaps

- Posterior inference allows us to discover the most salient gappy patterns in a corpus

  e.g., *not only __ but*     *either __ or*

- We validate the models by including patterns as features in a phrase-based MT system

- Code is available: `www.ark.cs.cmu.edu/MT`

# Motivation

- Gappy translation units have received a lot of attention recently
  - Mostly bilingual: Simard et al. (2005), Chiang (2005), Galley and Manning (2010), *inter alia*
  - But also monolingual: Xiong et al (2011)
- All rely on heuristics or mutual information
- Can we discover gappy patterns using generative models?

# Monolingual Pattern Models

- "Unigram" model: patterns generated independently

- Main intuition: words in a pattern are generated all at once

- Bayesian nonparametric priors and posterior inference favor the use of a small set of patterns to explain the data

nato must either say " yes " or " no " to the baltic states .

nato must either say " yes " or " no " to the baltic states .

- **Generative story:**
  - ☐ Generate number of word positions
  - ☐ Generate number of colors
  - ☐ Assign word positions to colors
  - ☐ Generate a lexical pattern for each color

nato must either say " yes " or " no " to the baltic states .

- **Generative story:**
  - ☐ Generate number of word positions (n = 16)
  - ☐ Generate number of colors
  - ☐ Assign word positions to colors
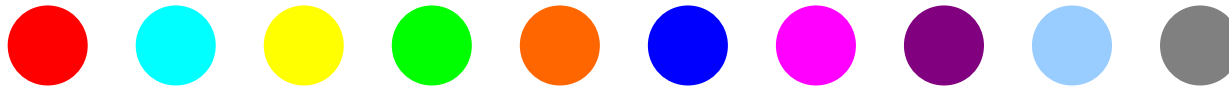  - ☐ Generate a lexical pattern for each color

nato must either say " yes " or " no " to the baltic states .

# Generative story:

- ☐ Generate number of word positions (n = 16)
- ☐ **Generate number of colors (m = 10)**
- ☐ Assign word positions to colors
- ☐ Generate a lexical pattern for each color

nato must either say " yes " or " no " to the baltic states .

# Generative story:

- ☐ Generate number of word positions (n = 16)
- ☐ Generate number of colors (m = 10)
- ☐ Assign word positions to colors
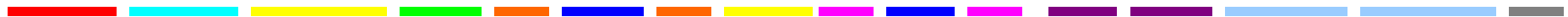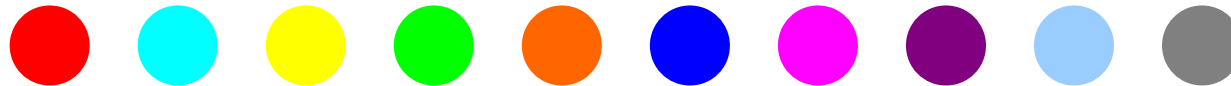- ☐ Generate a lexical pattern for each color

nato must either say " yes " or " no " to the baltic states .

# Generative story:

- Generate number of word positions (n = 16)
- Generate number of colors (m = 10)
- Assign word positions to colors
- Generate a lexical pattern for each color

nato must either say " yes " or " no " to the baltic states .

# ■ Generative story:

- ☐ Generate number of word positions (n = 16)

- ☐ Generate number of colors (m = 10)

- ☐ Assign word positions to colors
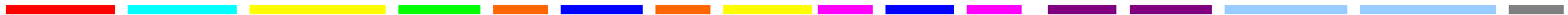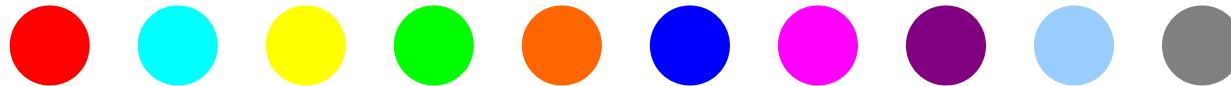
- ☐ **Generate a lexical pattern for each color**

nato

nato must either say " yes " or " no " to the baltic states .

# Generative story:

- Generate number of word positions (n = 16)
- Generate number of colors (m = 10)
- Assign word positions to colors
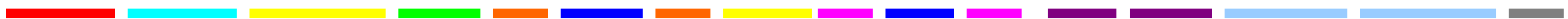- Generate a lexical pattern for each color

nato must

nato must either say " yes " or " no " to the baltic states .

# Generative story:

- Generate number of word positions (n = 16)
- Generate number of colors (m = 10)
- Assign word positions to colors
- Generate a lexical pattern for each color

nato must either                    or

# What is a **pattern**?

A sequence of symbols, possibly including the special symbol "__" which is used to indicate a gap of nonzero length

Examples:

nato                                        the united states

must                                        according to the __ ,

either __ or                            countries __ their __ the united states

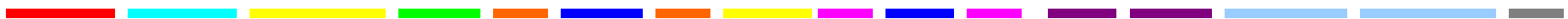nato must either                                    or

nato must either say " yes " or " no " to the baltic states .

# ■ Generative story:

- ☐ Generate number of word positions (n = 16)
- ☐ Generate number of colors (m = 10)
- ☐ Assign word positions to colors
- ☐ Generate a lexical pattern for each color
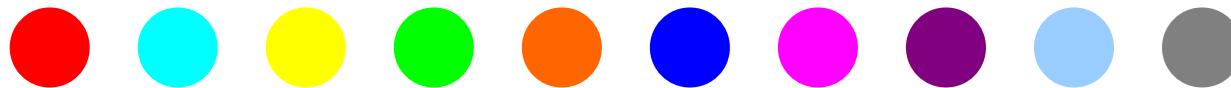
nato must either                          or                          baltic states

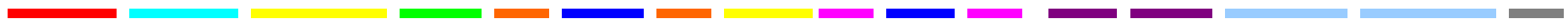nato must either say " yes " or " no " to the baltic states .

# ■ Generative story:

- ☐ Generate number of word positions (n = 16)
- ☐ Generate number of colors (m = 10)
- ☐ Assign word positions to colors
- ☐ Generate a lexical pattern for each color

nato must either          or  "        "          baltic states

nato must either say " yes " or " no " to the baltic states .

## ■ Generative story:

☐ Generate number of word positions (n = 16)

☐ Generate number of colors (m = 10)

☐ Assign word positions to colors

☐ Generate a lexical pattern for each color

nato must either say  "  yes  "  or  "  no  "  to  the baltic states  .
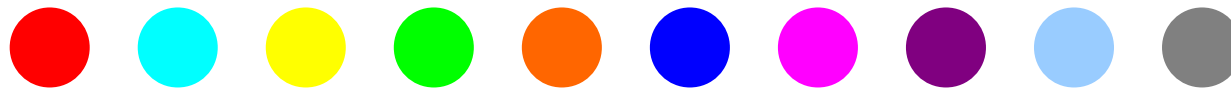
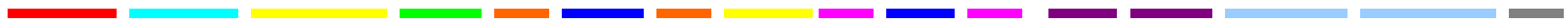nato must either say " yes " or " no " to the baltic states .

# Generative story:

- Generate number of word positions (n = 16)
- Generate number of colors (m = 10)
- Assign word positions to colors
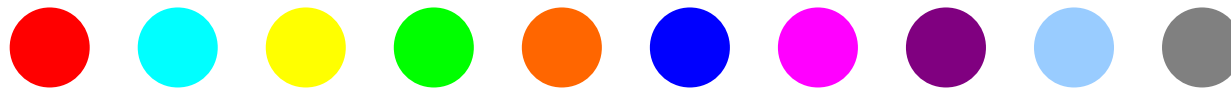- Generate a lexical pattern for each color

nato must either say " yes " or " no " to the baltic states .

nato must either say " yes " or " no " to the baltic states .

- **Generative story:**
  - ☐ Generate number of wo
  - ☐ Generate number of co
  - ☐ Assign word positions to colors
  - ☐ Generate a lexical pattern for each color

Uses a single multinomial distribution over patterns

nato must either say   "   yes   "   or   "   no   "   to  the baltic states   .

# Nonparametric Priors

- We use a single multinomial distribution over patterns ("unigram pattern model")
- Dirichlet process prior for this multinomial

# Inference

- ## Goal:
  - ☐ Given a corpus, obtain an estimate for how probable each pattern is

- ## To do this:
  - ☐ Obtain samples from posterior distribution over color assignments
  - ☐ Compute pattern counts from samples

# Inference

- **Goal:**
  - Given a corpus, obtain an estimate for how probable each pattern is

- **To do this:**
  - Obtain samples from posterior distribution over color assignments
  - Compute pattern counts from samples
  - We use collapsed Gibbs sampling to marginalize out the multinomial distribution over patterns

# Gibbs Sampling

- Go through each word and sample a new color

nato must either say " yes " or " no " to the baltic states .

# Gibbs Sampling

- Go through each word and sample a new color



nato must either say " yes " or " no " to the baltic states .

# Gibbs Sampling

- Go through each word and sample a new color
  - □ Choose any of the existing colors in the sentence, or
  - □ An entirely new color

nato must either say " yes " or " no " to the baltic states .

# Gibbs Sampling

- Go through each word and sample a new color
  - ☐ Choose any of the existing colors in the sentence, or
  - ☐ An entirely new color

nato must either say " yes " or " no " to the baltic states .

# Gibbs Sampling

■ We run sampling for 600 iterations on 125,000 sentences of English news commentary text

■ After burn-in, we average pattern counts across all samples

# Most Probable Patterns

| | | | |
|---|---|---|---|
| . | the __ " __ " | in __ , __ in | , however , |
| " __ " | the united states | america 's | " __ " __ " __ " |
| -- __ -- | rather than | more than | does not |
| ( __ ) | as __ as | china 's | why __ ? |
| the __ of | the __ of __ in | the __ of __ , | country 's |
| , __ , __ , | the __ is | prime minister | , __ the __ of |
| the __ ( __ ) | not only __ but | russia 's | its own |
| both __ and | the | europe 's | from __ to |
| have been | it is __ that | is | part of |
| the __ of __ and | should be | what __ ? | the __ between __ and |
| more __ than | their own | the world 's | such as __ , |
| - __ - | based on | between __ and | these |
| as well | of __ " __ " | developing countries | either __ or |
| this | not __ , but | climate change | economic growth |
| , __ " __ " | has been | the __ of __ 's | european union |

# Most Probable Patterns

| | | | |
|---|---|---|---|
| . | the __ " __ " | in __ , __ in | , however , |
| " __ " | the united states | america 's | " __ " __ " __ " |
| -- __ -- | rather than | more than | does not |
| ( __ ) | as __ as | china 's | why __ ? |
| the __ of | the __ of __ in | the __ of __ , | country 's |
| , __ , __ , | the __ is | prime minister | , __ the __ of |
| the __ ( __ ) | not only __ but | russia 's | its own |
| both __ and | the | europe 's | from __ to |
| have been | it is __ that | is | part of |
| the __ of __ and | should be | what __ ? | the __ between __ and |
| more __ than | their own | the world 's | such as __ , |
| - __ - | based on | between __ and | these |
| as well | of __ " __ " | developing countries | either __ or |
| this | not __ , but | climate change | economic growth |
| , __ " __ " | has been | the __ of __ 's | european union |

# Sorting by Conditional Probability

| | | |
|---|---|---|
| academy ___ sciences | treasury secretary ___ geithner | at ___ ghraib |
| beijing ___ shanghai | sooner ___ later | rule ___ law |
| booms ___ busts | first ___ foremost | free ___ fair |
| council ___ advisers | played ___ role | neither ___ nor |
| dominicans ___ haitian | down ___ road | across ___ border |
| flemish ___ walloons | freedom ___ expression | clash ___ civilizations |
| gref ___ program | at ___ disposal | estonia ___ lithuania |
| heat ___ droughts | take ___ granted | within ___ framework |
| humanitarian ___ displaced | - ___ - | window ___ opportunity |
| karnofsky ___ hassenfeld | at ___ expense | solve ___ problem |
| kazakhstan ___ kyrgyzstan | taken ___ granted | paid ___ price |
| portugal ___ greece | billions ___ dollars | taking ___ account |
| regulators ___ supervisors | answer ___ yes | during ___ period |
| sine ___ non | poland ___ slovakia | lender ___ last |
| stalin ___ mao | ukraine ___ orange | positive ___ negative |

# Using a Product of Experts

| | | | |
|---|---|---|---|
| -- __ -- | how __ ? | we __ our | his __ his |
| ( __ ) | the __ ( __ ) | over __ past | some __ others |
| - __ - | on __ basis | prevent __ from | may __ be |
| both __ and | less __ than | in __ way | as __ as |
| not only __ but | on __ other hand | one __ another | oil __ gas |
| " __ " | at __ level | political __ economic | at __ moment |
| more __ than | it is __ that | for __ reasons | such as __ and |
| either __ or | not __ , but | at __ time | question __ whether |
| why __ ? | play __ role | more __ more | if __ then |
| neither __ nor | france __ germany | the rest __ world | war __ iraq |
| what __ ? | he __ his | more __ less | ; __ ; |
| rule __ law | allow __ to | in __ region | have __ been |
| whether __ or | for __ first time | rich __ poor | in __ cases |
| around __ world | china __ india | as __ whole | war __ terror |
| has __ been | what __ do | on __ scale | at __ cost |

# Punctuation

| | | | |
|---|---|---|---|
| -- __ -- | how __ ? | we __ our | his __ his |
| ( __ ) | the __ ( __ ) | over __ past | some __ others |
| - __ - | on __ basis | prevent __ from | may __ be |
| both __ and | less __ than | in __ way | as __ as |
| not only __ but | on __ other hand | one __ another | oil __ gas |
| " __ " | at __ level | political __ economic | at __ moment |
| more __ than | it is __ that | for __ reasons | such as __ and |
| either __ or | not __ , but | at __ time | question __ whether |
| why __ ? | play __ role | more __ more | if __ then |
| neither __ nor | france __ germany | the rest __ world | war __ iraq |
| what __ ? | he __ his | more __ less | ; __ ; |
| rule __ law | allow __ to | in __ region | have __ been |
| whether __ or | for __ first time | rich __ poor | in __ cases |
| around __ world | china __ india | as __ whole | war __ terror |
| has __ been | what __ do | on __ scale | at __ cost |

# Connectives and Constructions

| | | | |
|---|---|---|---|
| -- __ -- | how __ ? | we __ our | his __ his |
| ( __ ) | the __ ( __ ) | over __ past | some __ others |
| - __ - | on __ basis | prevent __ from | may __ be |
| both __ and | less __ than | in __ way | as __ as |
| not only __ but | on __ other hand | one __ another | oil __ gas |
| " __ " | at __ level | political __ economic | at __ moment |
| more __ than | it is __ that | for __ reasons | such as __ and |
| either __ or | not __ , but | at __ time | question __ whether |
| why __ ? | play __ role | more __ more | if __ then |
| neither __ nor | france __ germany | the rest __ world | war __ iraq |
| what __ ? | he __ his | more __ less | ; __ ; |
| rule __ law | allow __ to | in __ region | have __ been |
| whether __ or | for __ first time | rich __ poor | in __ cases |
| around __ world | china __ india | as __ whole | war __ terror |
| has __ been | what __ do | on __ scale | at __ cost |

# Agreement

| | | | |
|---|---|---|---|
| -- __ -- | how __ ? | we __ our | his __ his |
| ( __ ) | the __ ( __ ) | over __ past | some __ others |
| - __ - | on __ basis | prevent __ from | may __ be |
| both __ and | less __ than | in __ way | as __ as |
| not only __ but | on __ other hand | one __ another | oil __ gas |
| " __ " | at __ level | political __ economic | at __ moment |
| more __ than | it is __ that | for __ reasons | such as __ and |
| either __ or | not __ , but | at __ time | question __ whether |
| why __ ? | play __ role | more __ more | if __ then |
| neither __ nor | france __ germany | the rest __ world | war __ iraq |
| what __ ? | he __ his | more __ less | ; __ ; |
| rule __ law | allow __ to | in __ region | have __ been |
| whether __ or | for __ first time | rich __ poor | in __ cases |
| around __ world | china __ india | as __ whole | war __ terror |
| has __ been | what __ do | on __ scale | at __ cost |

# Topicality

| | | | |
|---|---|---|---|
| -- __ -- | how __ ? | we __ our | his __ his |
| ( __ ) | the __ ( __ ) | over __ past | some __ others |
| - __ - | on __ basis | prevent __ from | may __ be |
| both __ and | less __ than | in __ way | as __ as |
| not only __ but | on __ other hand | one __ another | oil __ gas |
| " __ " | at __ level | political __ economic | at __ moment |
| more __ than | it is __ that | for __ reasons | such as __ and |
| either __ or | not __ , but | at __ time | question __ whether |
| why __ ? | play __ role | more __ more | if __ then |
| neither __ nor | france __ germany | the rest __ world | war __ iraq |
| what __ ? | he __ his | more __ less | ; __ ; |
| rule __ law | allow __ to | in __ region | have __ been |
| whether __ or | for __ first time | rich __ poor | in __ cases |
| around __ world | china __ india | as __ whole | war __ terror |
| has __ been | what __ do | on __ scale | at __ cost |

# Prepositional Phrases

| | | | |
|---|---|---|---|
| -- ___ -- | how ___ ? | we ___ our | his ___ his |
| ( ___ ) | the ___ ( ___ ) | over ___ past | some ___ others |
| - ___ - | on ___ basis | prevent ___ from | may ___ be |
| both ___ and | less ___ than | in ___ way | as ___ as |
| not only ___ but | on ___ other hand | one ___ another | oil ___ gas |
| " ___ " | at ___ level | political ___ economic | at ___ moment |
| more ___ than | it is ___ that | for ___ reasons | such as ___ and |
| either ___ or | not ___ , but | at ___ time | question ___ whether |
| why ___ ? | play ___ role | more ___ more | if ___ then |
| neither ___ nor | france ___ germany | the rest ___ world | war ___ iraq |
| what ___ ? | he ___ his | more ___ less | ; ___ ; |
| rule ___ law | allow ___ to | in ___ region | have ___ been |
| whether ___ or | for ___ first time | rich ___ poor | in ___ cases |
| around ___ world | china ___ india | as ___ whole | war ___ terror |
| has ___ been | what ___ do | on ___ scale | at ___ cost |

- How does this differ from word trigger pairs derived from mutual information (Rosenfeld, 1994)?
- The X ___ Y pairs we extract are similar to his pairs
- We also model collocations and larger patterns:
  - X Y Z
  - X Y ___ Z
  - X ___ Y ___ Z
  - X Y ___ Z ___ W ___ V
  - etc.
- Generative models are also amenable to extensions…

# Modeling Bilingual Patterns

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

# Modeling Bilingual Patterns

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- **Generative story:**
  - ☐ Generate n target word positions, n′ source positions
  - ☐ Generate m target colors, m′ source-only colors
  - ☐ Generate 1-to-1 word alignment between the word positions
  - ☐ Assign target word positions to target colors
  - ☐ Assign source word positions to either source colors or target colors
  - ☐ Generate a lexical pattern for each color

**Carnegie Mellon**

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- Generative story:
  - Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - Generate m target colors, m′ source-only colors
  - Generate 1-to-1 word alignment between the word positions
  - Assign target word positions to target colors
  - Assign source word positions to either source colors or target colors
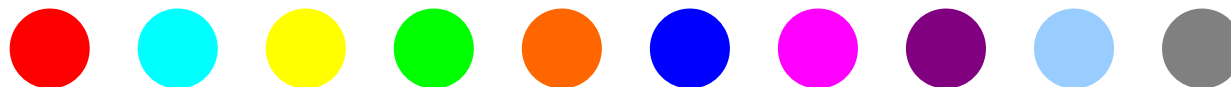  - Generate a lexical pattern for each color

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- **Generative story:**
  - ☐ Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - ☐ **Generate m target colors (m = 10), m′ source-only colors (m′ = 4)**
  - ☐ Generate 1-to-1 word alignment between the word positions
  - ☐ Assign target word positions to target colors
  - ☐ Assign source word positions to either source colors or target colors
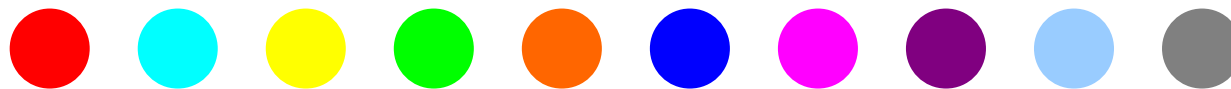  - ☐ Generate a lexical pattern for each color

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- Generative story:
  - Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
  - **Generate 1-to-1 word alignment between the word positions**
  - Assign target word positions to target colors
  - Assign source word positions to either source colors or target colors
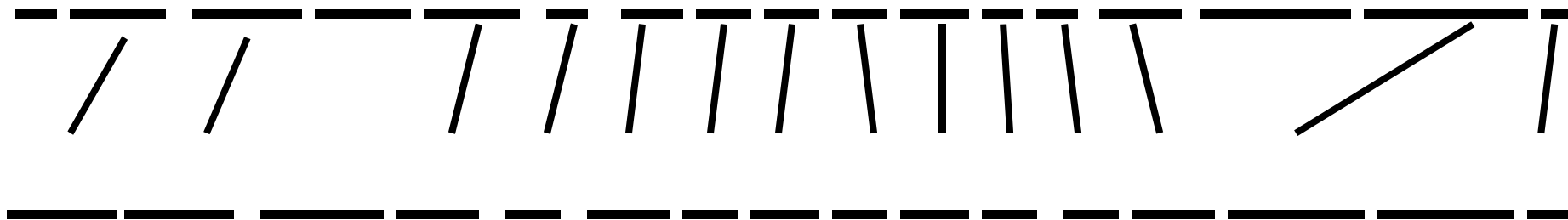  - Generate a lexical pattern for each color

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- Generative story:
  - ☐ Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - ☐ Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
  - ☐ Generate 1-to-1 word alignment between the word positions
  - ☐ **Assign target word positions to target colors**
  - ☐ Assign source word positions to either source colors or target colors
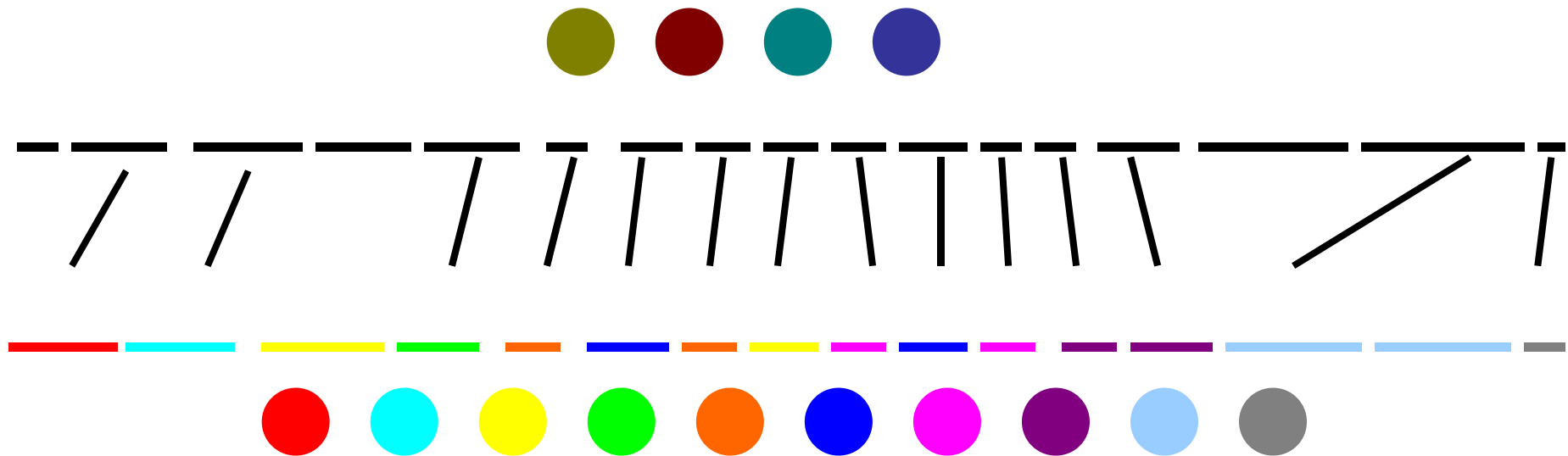  - ☐ Generate a lexical pattern for each color

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- ■ Generative story:
    - □ Generate n target word positions (n = 16), n′ source positions (n′ = 17)
    - □ Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
    - □ Generate 1-to-1 word alignment between the word positions
    - □ Assign target word positions to target colors
    - □ **Assign source word positions to either source colors or target colors**
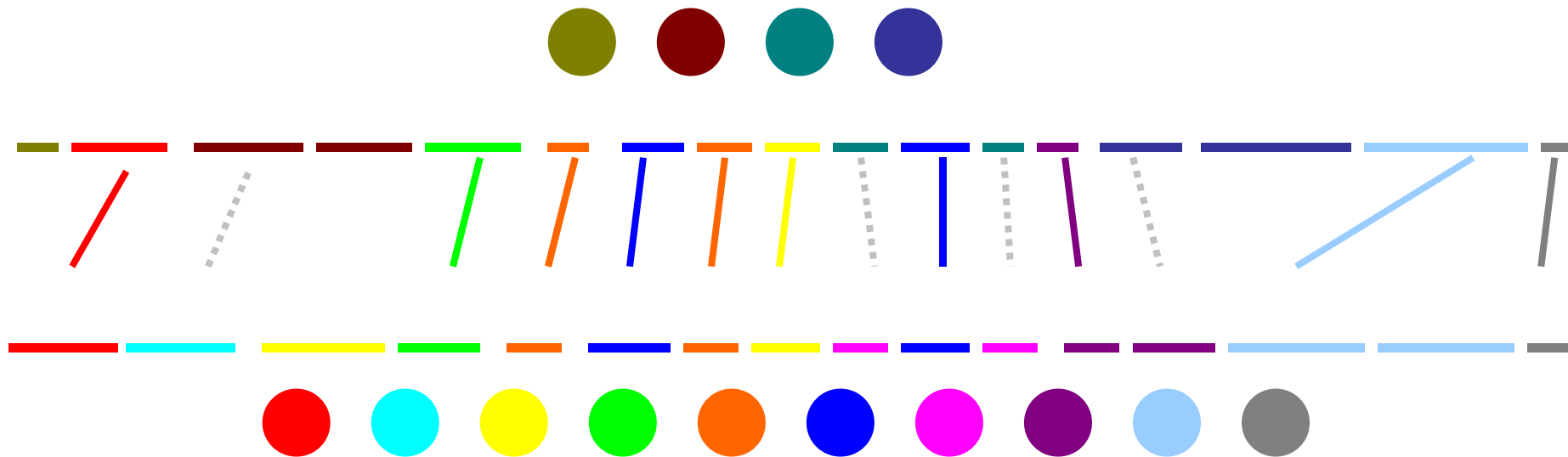    - □ Generate a lexical pattern for each color

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- ■ Generative story:
    - □ Generate n target word positions (n = 16), n′ source positions (n′ = 17)
    - □ Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
    - □ Generate 1-to-1 word alignment between the word positions
    - □ Assign target word positions to target colors
    - □ Assign source word positions to either source colors or target colors
    - □ **Generate a lexical pattern for each color**
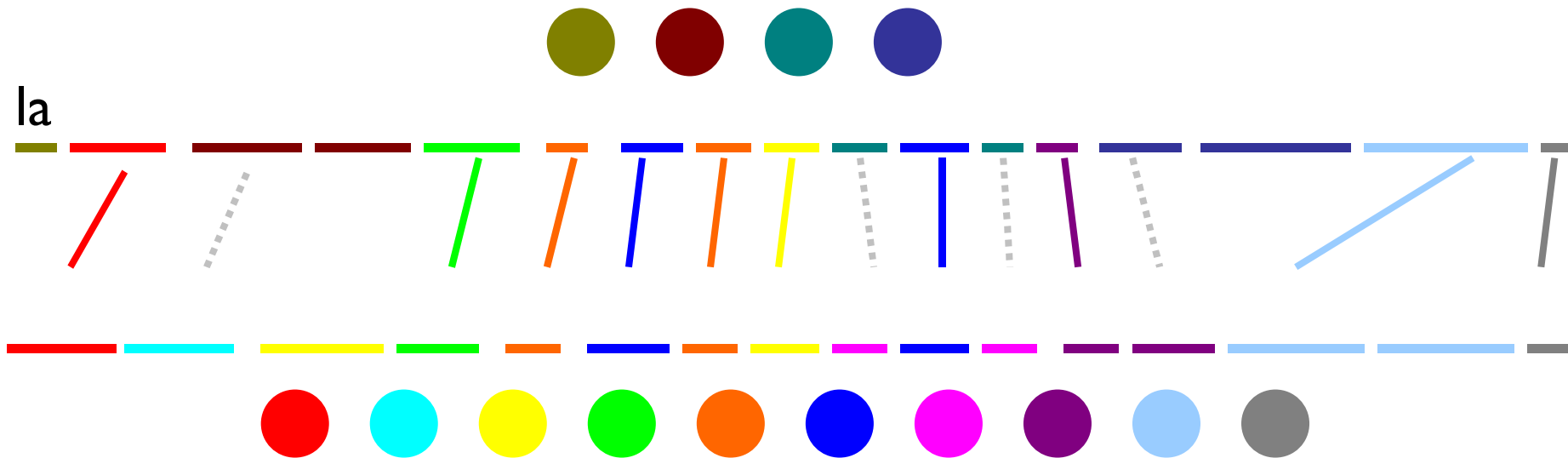
la

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- ■ Generative story:
  - □ Generate n target word positions (n = 16), n' source positions (n' = 17)
  - □ Generate m target colors (m = 10), m' source-only colors (m' = 4)
  - □ Generate 1-to-1 word alignment between the word positions
  - □ Assign target word positions to target colors
  - □ Assign source word positions to either source colors or target colors
  - □ **Generate a lexical pattern for each color**
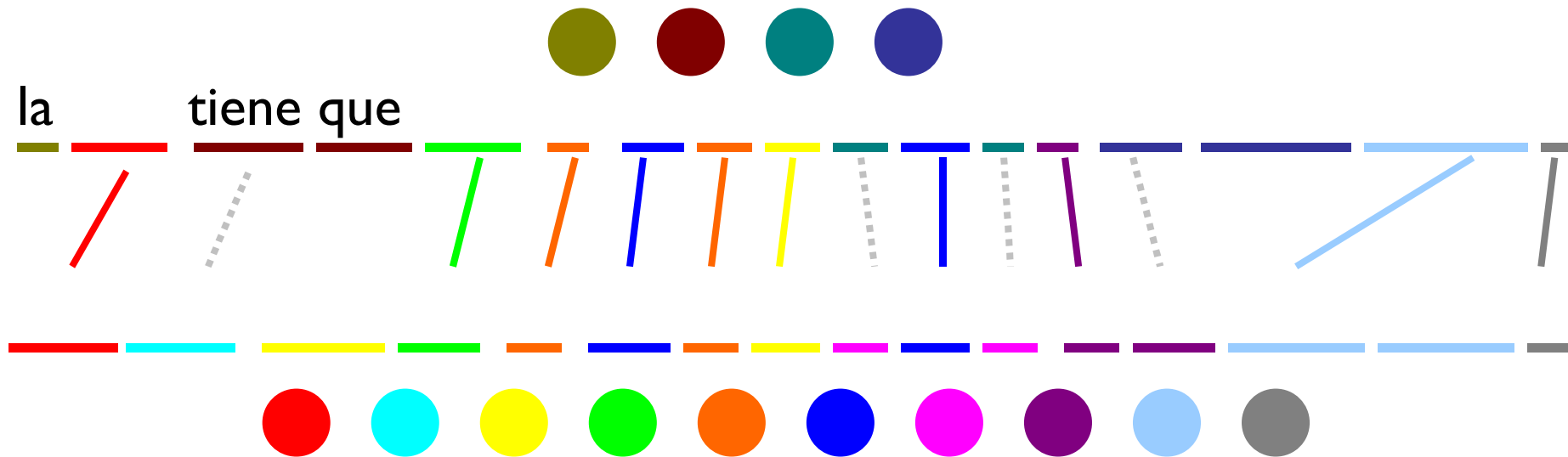
la          tiene que

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- ■ Generative story:
  - □ Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - □ Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
  - □ Generate 1-to-1 word alignment between the word positions
  - □ Assign target word positions to target colors
  - □ Assign source word positions to either source colors or target colors
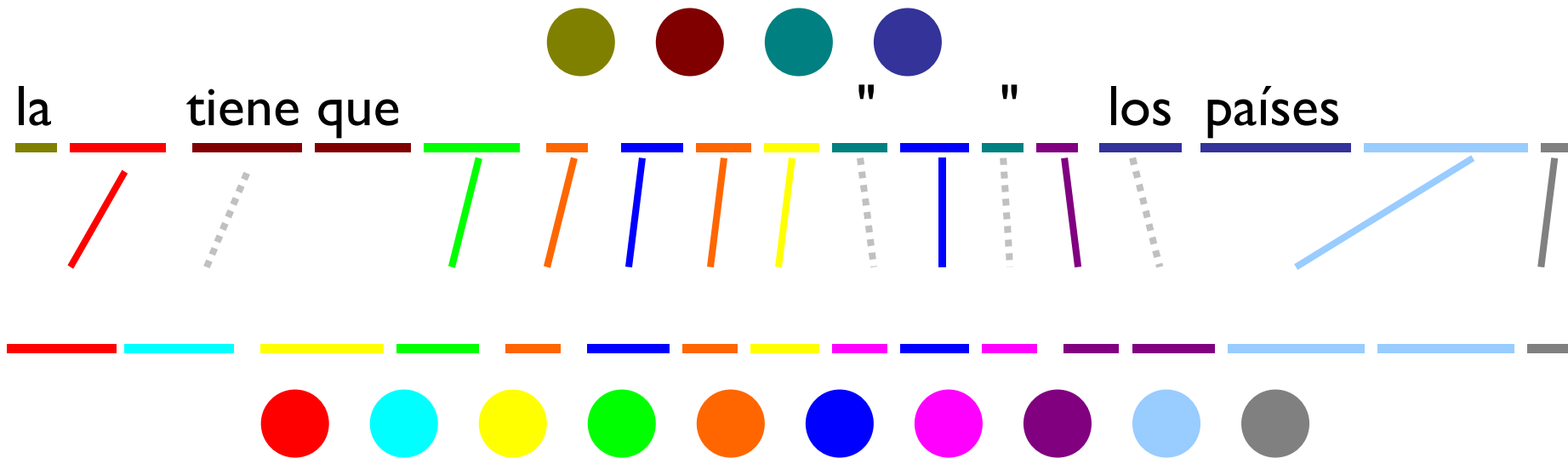  - □ Generate a lexical pattern for each color

la     tiene que     "     "     los países

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- **Generative story:**
  - Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
  - Generate 1-to-1 word alignment between the word positions
  - Assign target word positions to target colors
  - Assign source word positions to either source colors or target colors
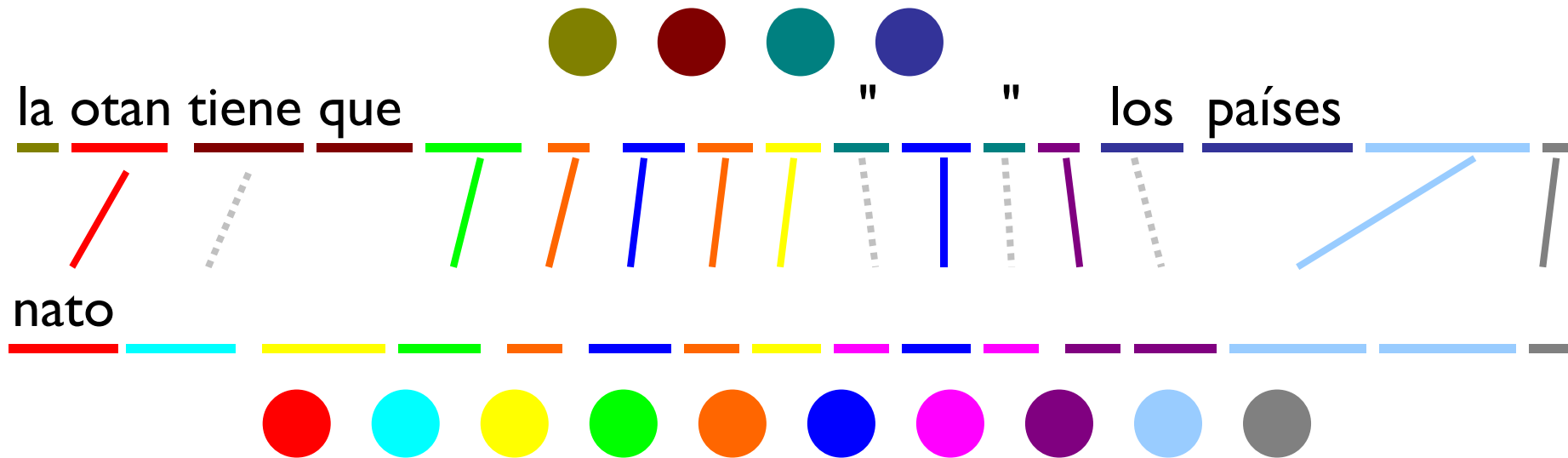  - Generate a lexical pattern for each color

la otan tiene que           "     "     los  países

nato

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- Generative story:
  - Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
  - Generate 1-to-1 word alignment between the word positions
  - Assign target word positions to target colors
  - Assign source word positions to either source colors or target colors
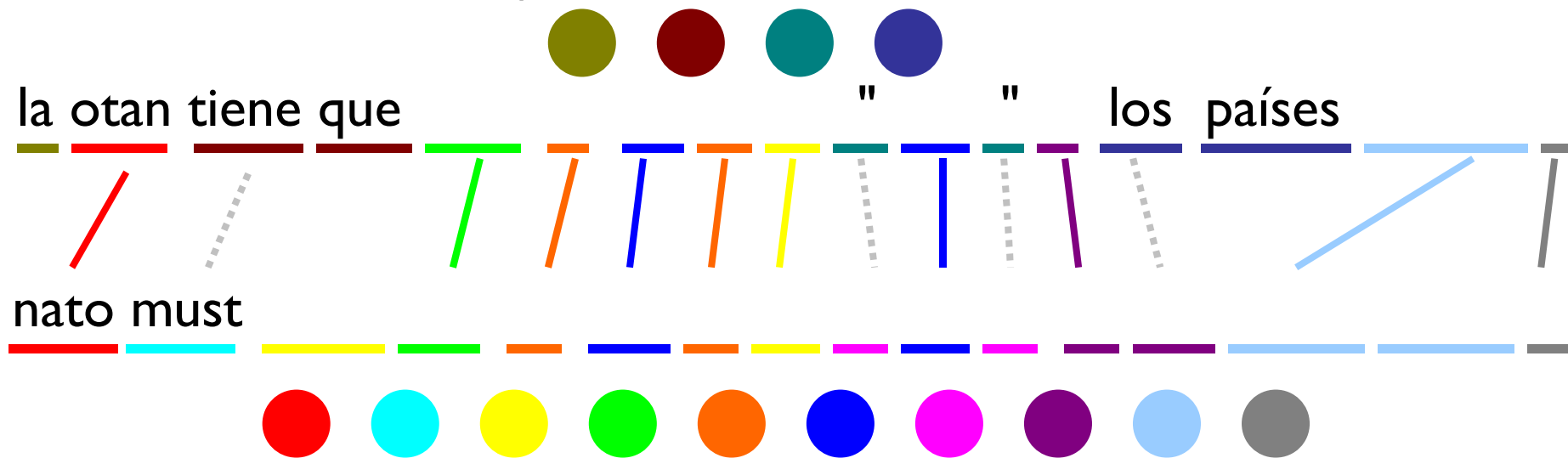  - Generate a lexical pattern for each color

la otan tiene que        "       "        los  países

nato must

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- ■ Generative story:
  - □ Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - □ Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
  - □ Generate 1-to-1 word alignment between the word positions
  - □ Assign target word positions to target colors
  - □ Assign source word positions to either source colors or target colors
  - □ Generate a lexical pattern for each color

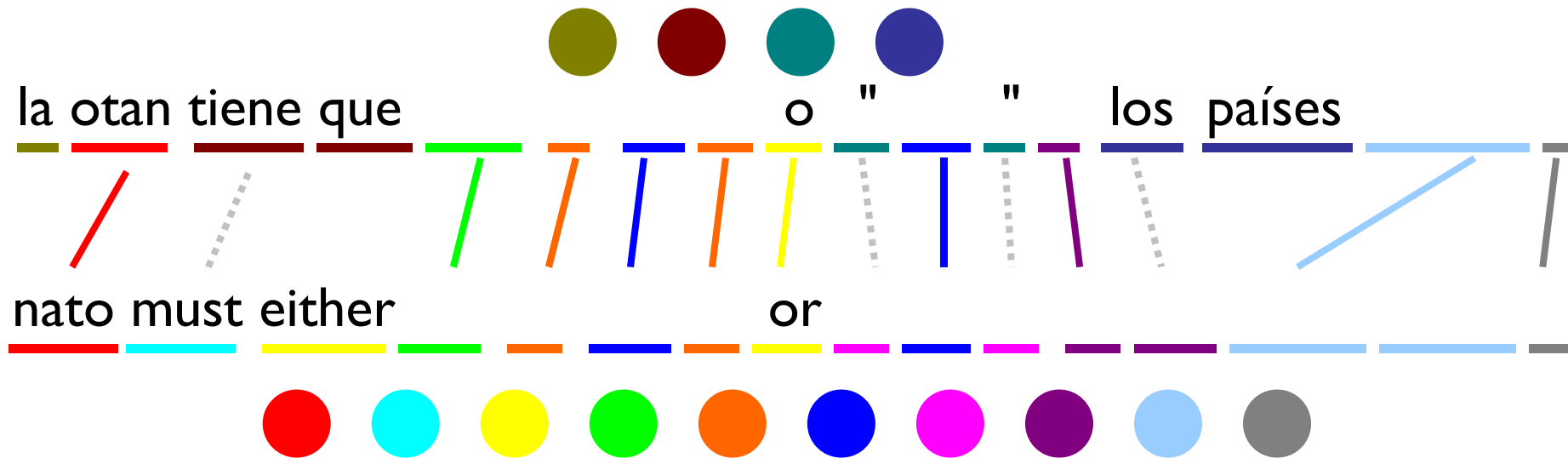la otan tiene que        o "   "   los países

nato must either       or

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- Generative story:
  - Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
  - Generate 1-to-1 word alignment between the word positions
  - Assign target word positions to target colors
  - Assign source word positions to either source colors or target colors
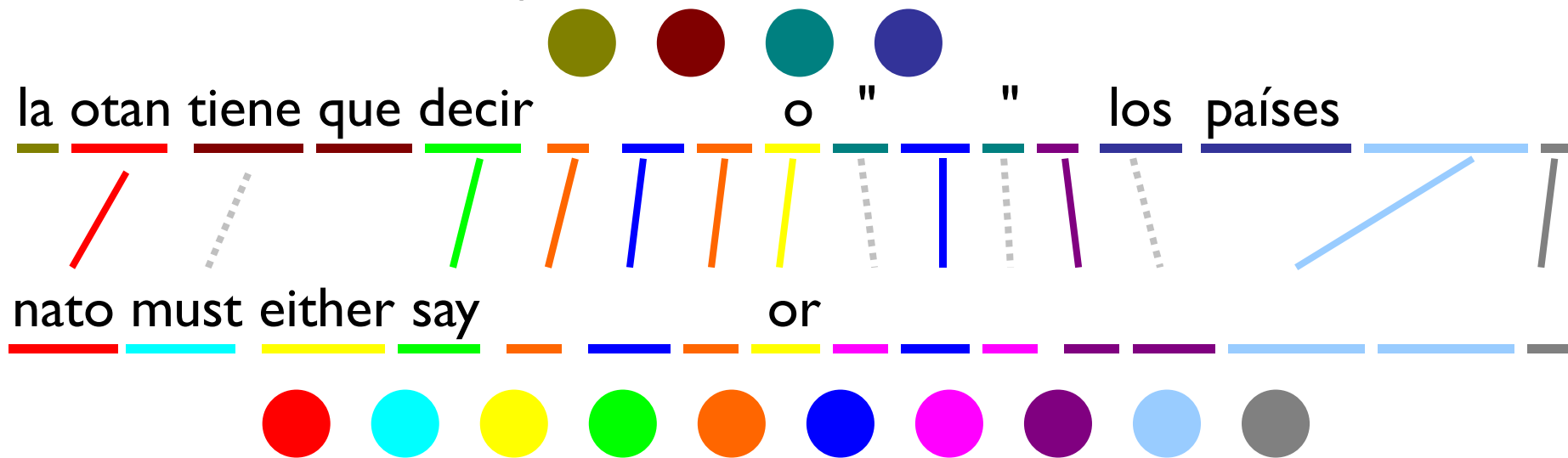  - Generate a lexical pattern for each color

la otan tiene que decir            o  "        "    los  países
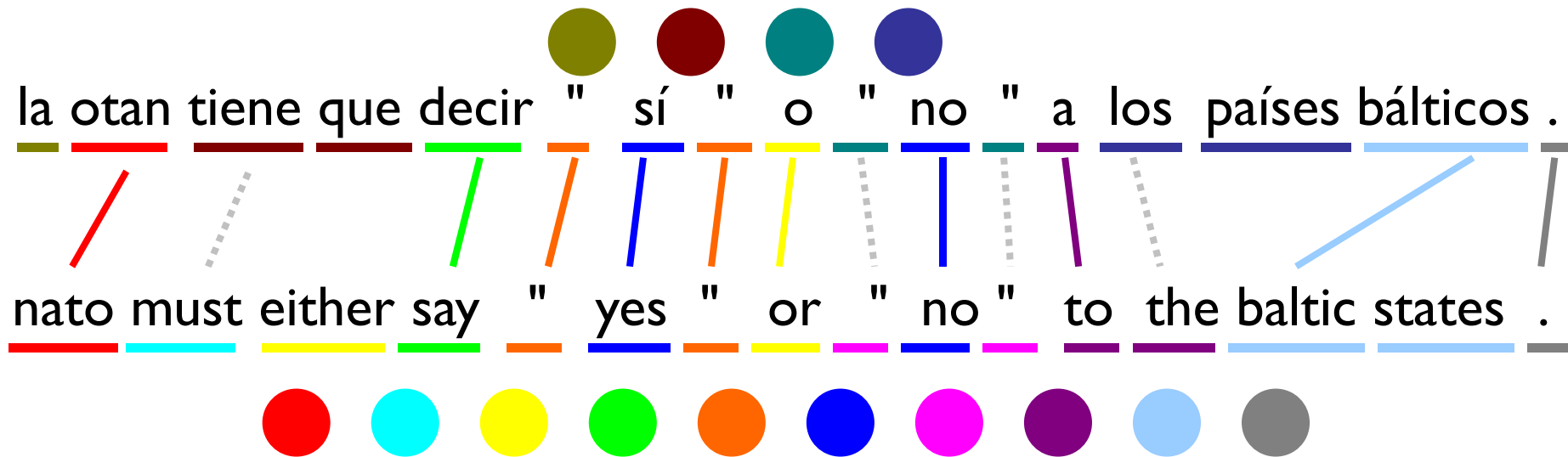
nato must either say            or

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

- ■ Generative story:
  - □ Generate n target word positions (n = 16), n′ source positions (n′ = 17)
  - □ Generate m target colors (m = 10), m′ source-only colors (m′ = 4)
  - □ Generate 1-to-1 word alignment between the word positions
  - □ Assign target word positions to target colors
  - □ Assign source word positions to either source colors or target colors
  - □ Generate a lexical pattern for each color

la otan tiene que decir " sí " o " no " a los países bálticos .

nato must either say " yes " or " no " to the baltic states .

# Inference with Bilingual Pattern Models

- Gibbs sampler similar to monolingual model, with a few extra moves (see paper and code)

- Inference run for 300 iterations

- Examples:

| | |
|---|---|
| we must(debemos) | they ___ their(sus) |
| we are(estamos) | their(sus) ___ their(sus) |
| we can(podemos) | he ___ his(sus) |
| either ___ or(o) | it ___ its(sus) |

# Experiments

- We added count features for individual patterns

- Monolingual patterns:
  - 10k lexical patterns
  - 15k patterns on Brown clusters

- Bilingual patterns:
  - 5k word/word, 5k word/cluster, 5k cluster/cluster

- Features are **non-local** since they can match anywhere in the derivation

- Features incorporated via cube pruning of phrase lattices

- Trained using a MIRA-like procedure (simplified from that of Chiang et al., 2009)

# Experiments

- **Chinese-English**
  - ☐ 300k sentence pairs from FBIS corpus
  - ☐ Tuned on MT03, tested on MT05
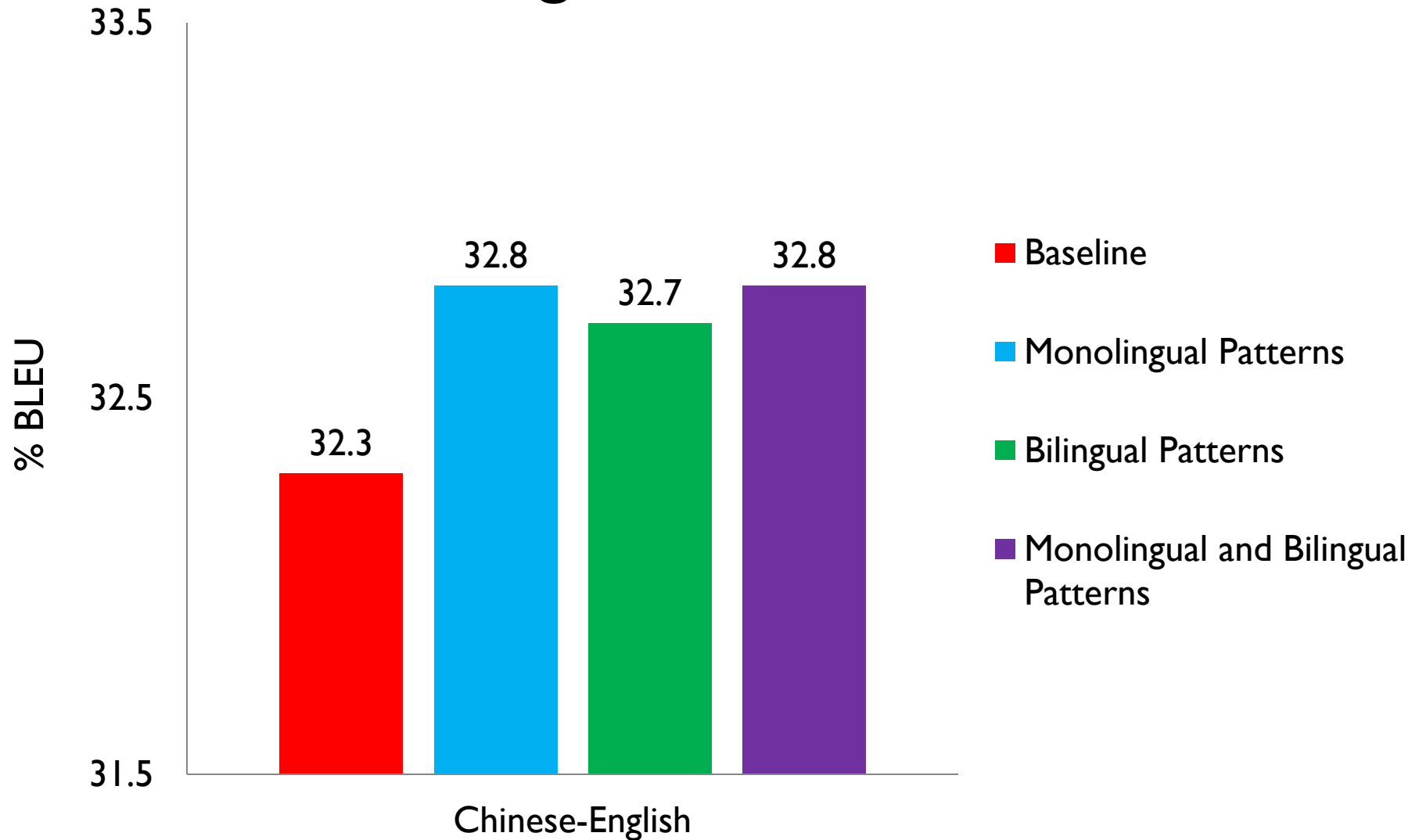  - ☐ Trigram LM estimated from English side of parallel corpus + 200M words of Gigaword data

- **Spanish-English**
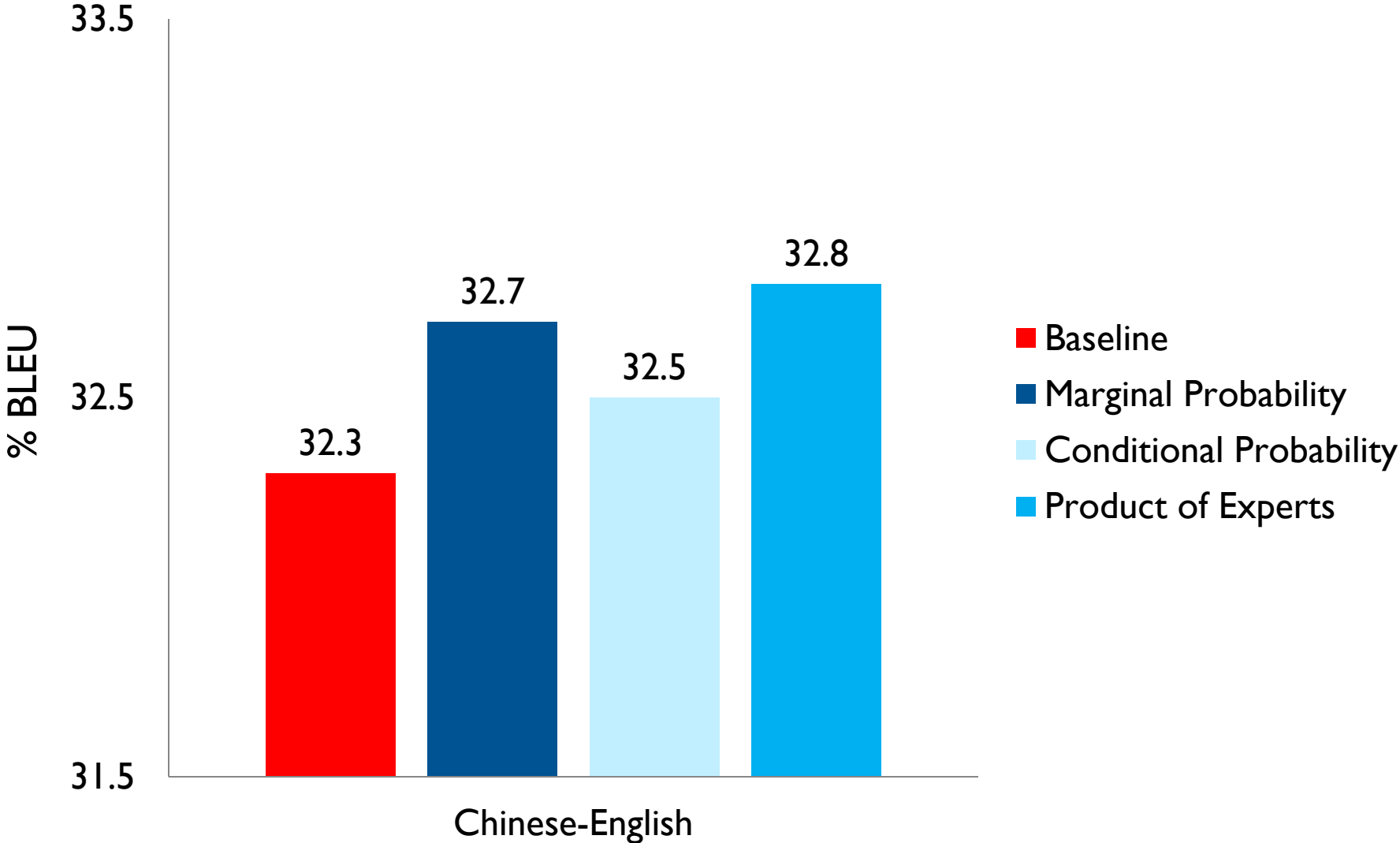  - ☐ No improvement; experiments reported in paper

# Adding Pattern Features



Chart title: Adding Pattern Features

Y-axis: % BLEU (scale from 31.5 to 33.5)

X-axis: Chinese-English

Bar values:
- Baseline (red): 32.3
- Monolingual Patterns (blue): 32.8
- Bilingual Patterns (green): 32.7
- Monolingual and Bilingual Patterns (purple): 32.8

Comparing Ways of Ranking Patterns

# Most Highly-Weighted Features

| | |
|---|---|
| said that ___ the | of ___ million |
| however , ___ the | , ___ likely |
| agence france ___ presse | said that ___ and |
| 's ___ , ___ 's | added ___ " |
| us ___ iraq | - ___ - |
| reported ___ the | rate ___ percent |

the ___ {media, school, university, election, bank} ___ {made, established, given, taken, reached}

{said, stressed, stated, indicated, noted} that ___ in

{meeting, report, conference, reports, summit} ___ {1, july, june, march, april}

{news, press, spokesman, reporter, consultative} {meeting, ...} ___ {1, july, june, march, april}

{news, press, spokesman, reporter, consultative} ___ {1, july, june, march, april}

the ___ {enterprises, companies, students, customers, others} ___ {enterprises, companies, ...}

{japan, russia, europe, 2003, 2004} ___ {us, japanese, russian, u.s., british}

# Conclusions

- We presented models for discovering gappy patterns in monolingual and parallel text

- Validation of patterns qualitatively and quantitatively in a phrase-based MT system

- Code implementing inference for our models is available: `www.ark.cs.cmu.edu/MT`

# Thanks!