

Word Alignment

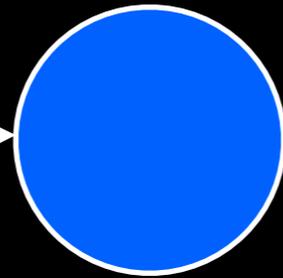
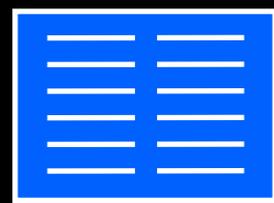
Adam Lopez
Johns Hopkins

Quick Recap

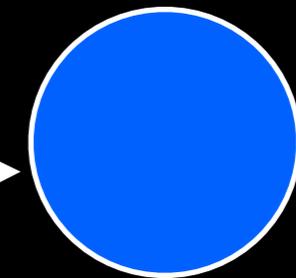
training data
(parallel text)

learner

TM + LM model



联合国安全理事会的
五个常任理事国都



decoder

However, the sky remained clear
under the strong north wind.

IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。

IBM Model 1

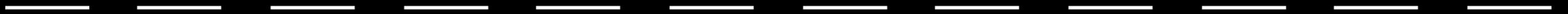
Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε



IBM Model 1

Although north wind howls , but sky still very clear .

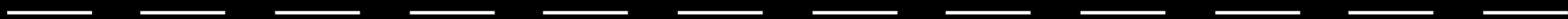
虽然北风呼啸，但天空依然十分清澈。ε

$$p(\text{English length} | \text{Chinese length})$$

IBM Model 1

Although north wind howls , but sky still very clear .

虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。 ϵ

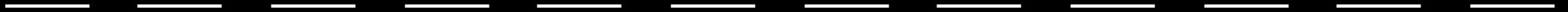


$$p(\text{English length} | \text{Chinese length})$$

IBM Model 1

Although north wind howls , but sky still very clear .

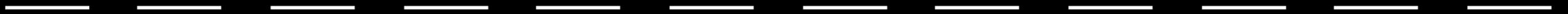
虽然北风呼啸，但天空依然十分清澈。ε



IBM Model 1

Although north wind howls , but sky still very clear .

虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。 ϵ

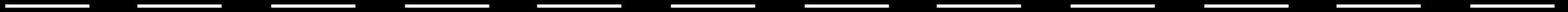


p(Chinese word position)

IBM Model 1

Although north wind howls , but sky still very clear .

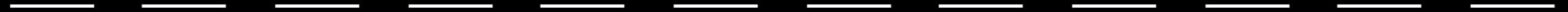
虽然北风呼啸，但天空依然十分清澈。ε



IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε



However

IBM Model 1

Although north wind howls , but sky still very clear .

虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。 ϵ



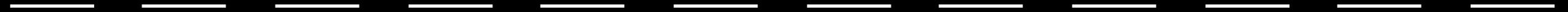
However

$$p(\textit{English word} | \textit{Chinese word})$$

IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

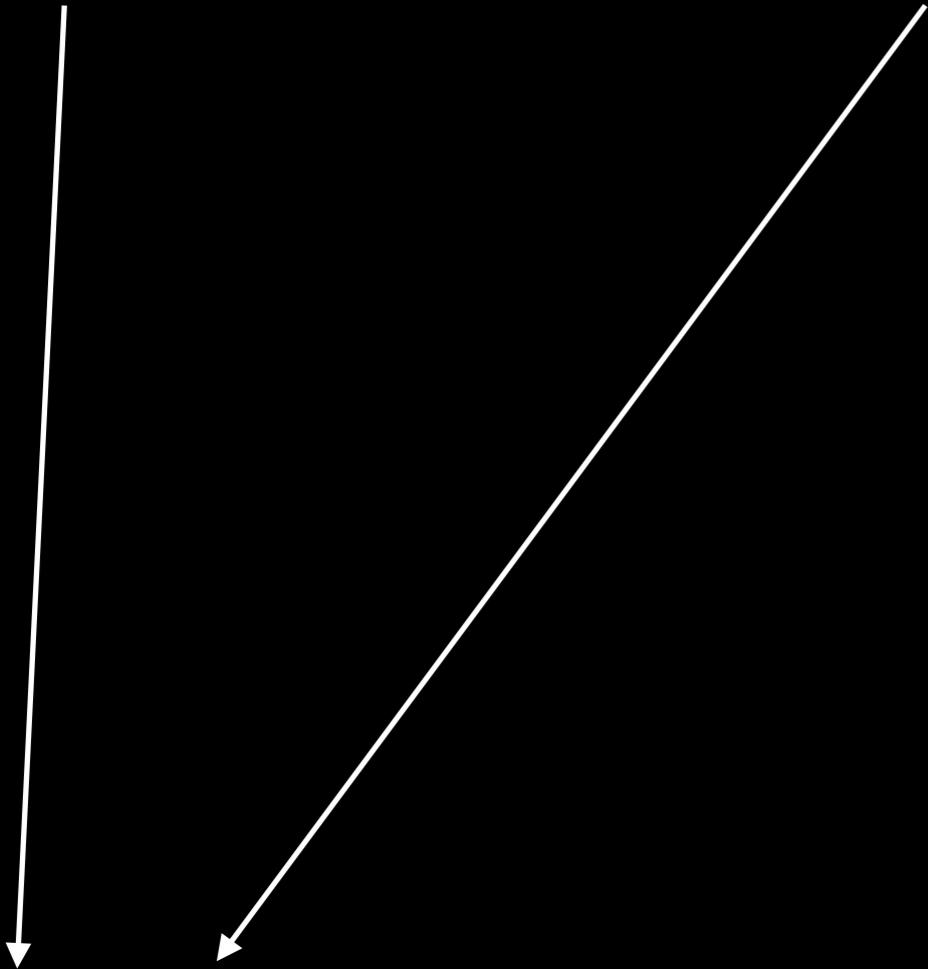


However

IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

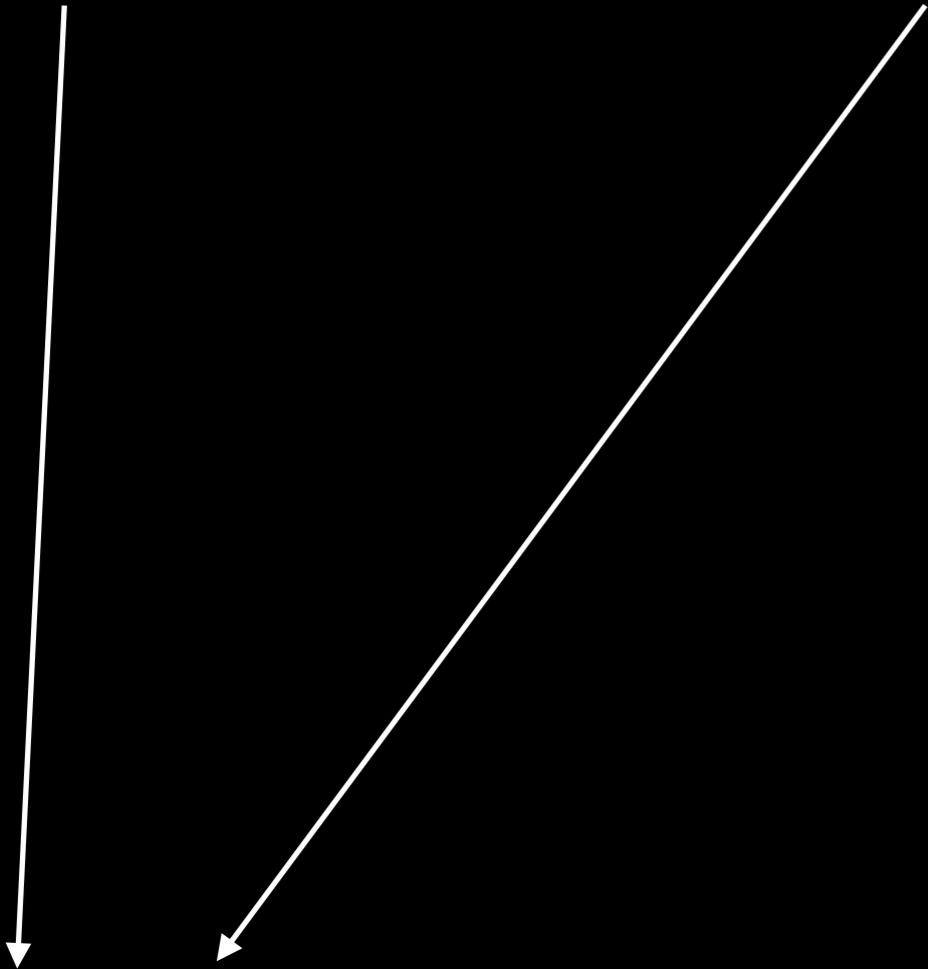


However

IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

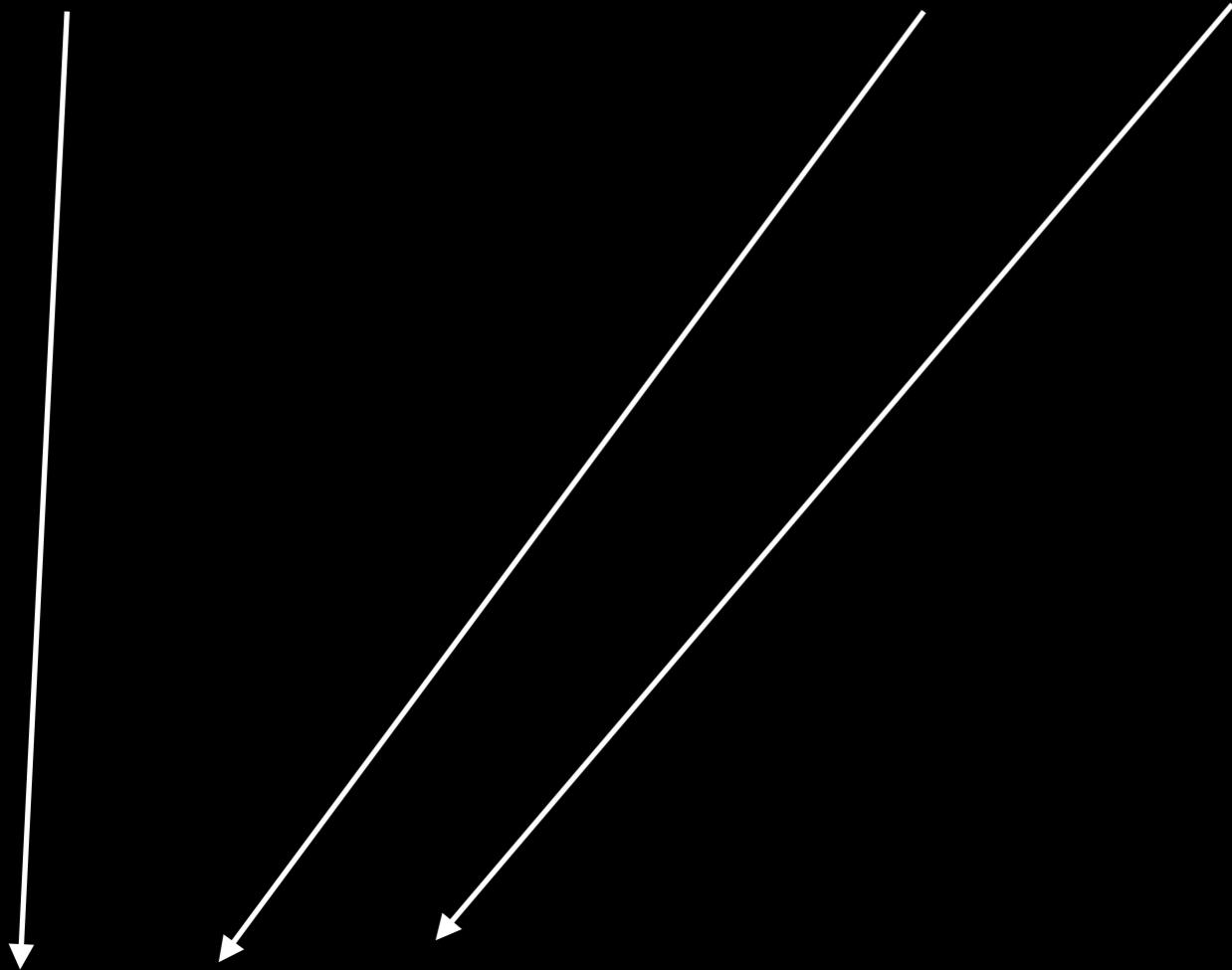


However ,

IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

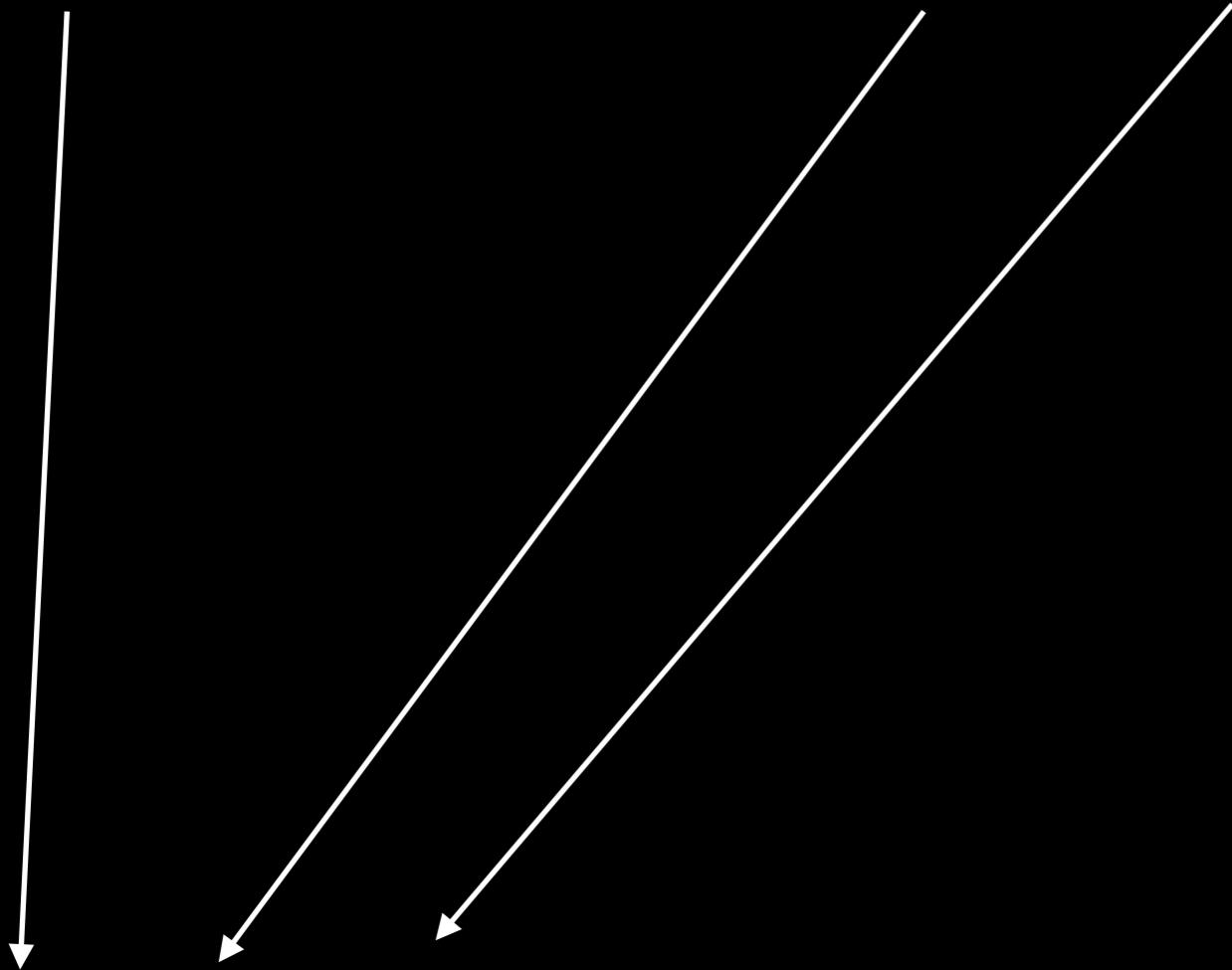


However ,

IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

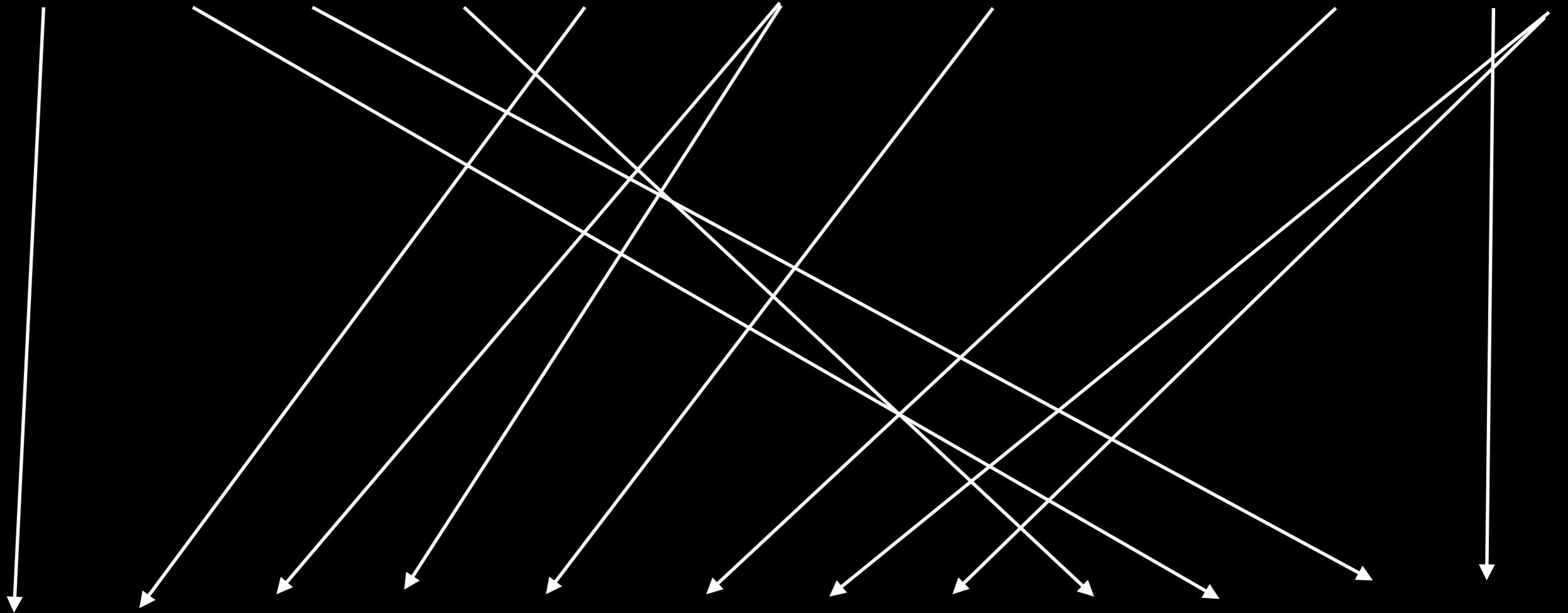


However , the

IBM Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε



However , the sky remained clear under the strong north wind .

IBM Model 1

$p(\textit{despite} | \text{虽然})$

$p(\textit{however} | \text{虽然})$

$p(\textit{although} | \text{虽然})$

$p(\textit{northern} | \text{北})$

$p(\textit{north} | \text{北})$

IBM Model 1

$p(\textit{despite} | \text{虽然})$???

$p(\textit{however} | \text{虽然})$???

$p(\textit{although} | \text{虽然})$???

$p(\textit{northern} | \text{北})$???

$p(\textit{north} | \text{北})$???

IBM Model 1

θ	{	$p(\textit{despite} \text{虽然})$???
		$p(\textit{however} \text{虽然})$???
		$p(\textit{although} \text{虽然})$???
		$p(\textit{northern} \text{北})$???
		$p(\textit{north} \text{北})$???





$p(\text{heads}) ?$

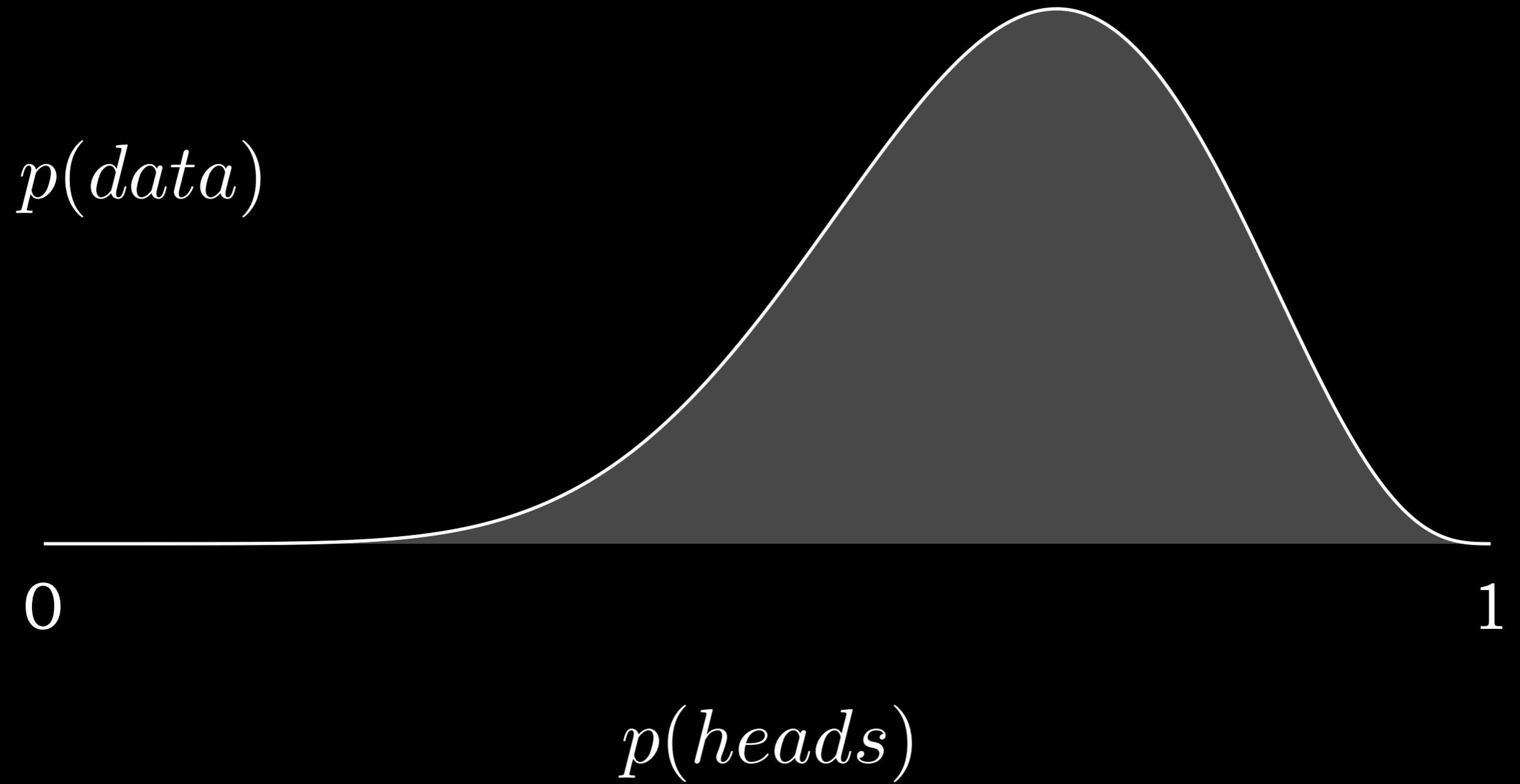


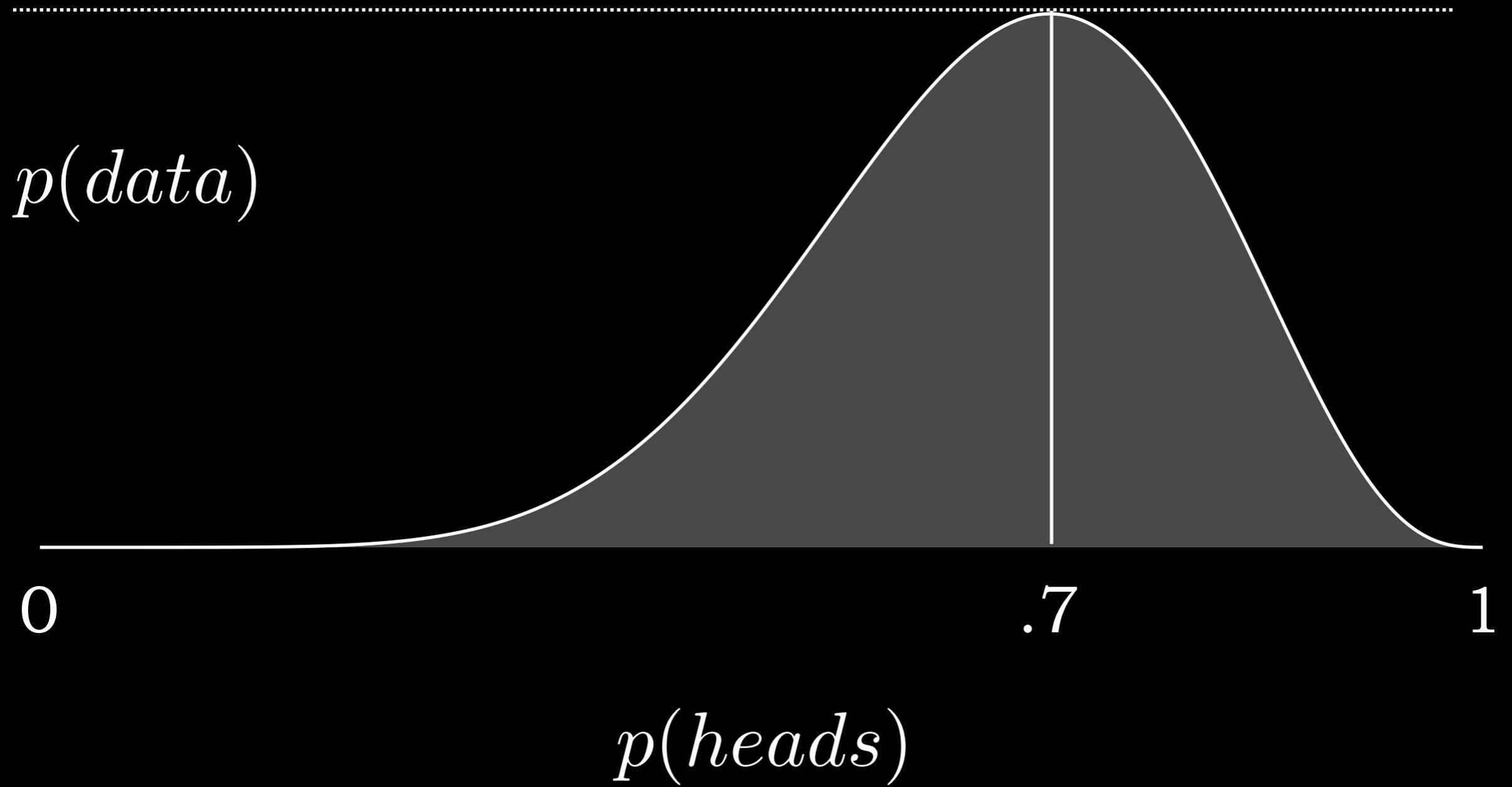


$$p(\text{data}) = p(\text{heads})^7 \times p(\text{tails})^3$$



$$p(\text{data}) = p(\text{heads})^7 \times [1 - p(\text{heads})]^3$$

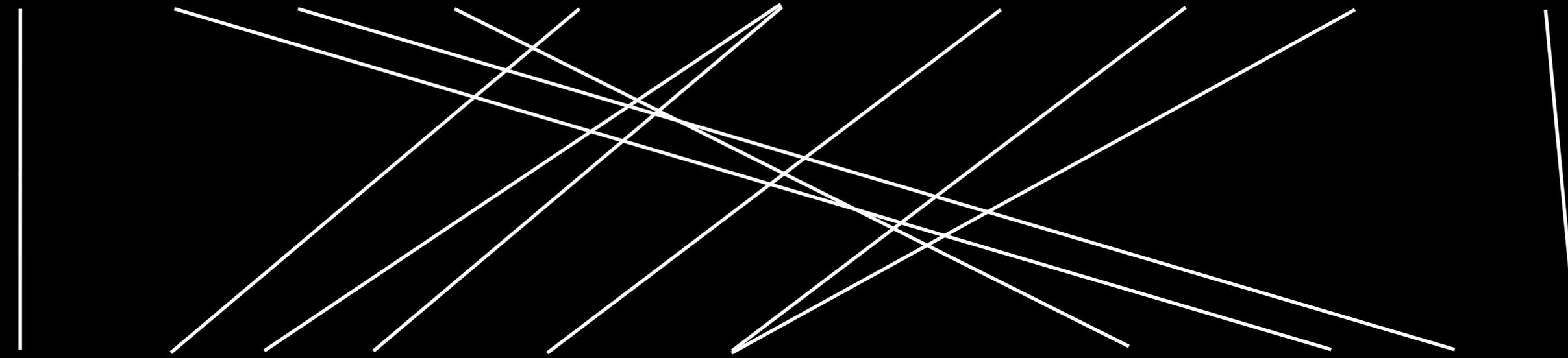




IBM Model 1

Although north wind howls , but sky still very clear .

虽然 北 风 呼啸 ， 但 天空 依然 十分 清澈 。



However , the sky remained clear under the strong north wind .

$$p(\textit{however} | \text{虽然}) = \frac{\text{\# of times 虽然 aligns to However}}{\text{\# of times 虽然 occurs}}$$

IBM Model 1

Although north wind howls , but sky still very clear .

虽然 北 风 呼啸 ， 但 天空 依然 十分 清澈 。

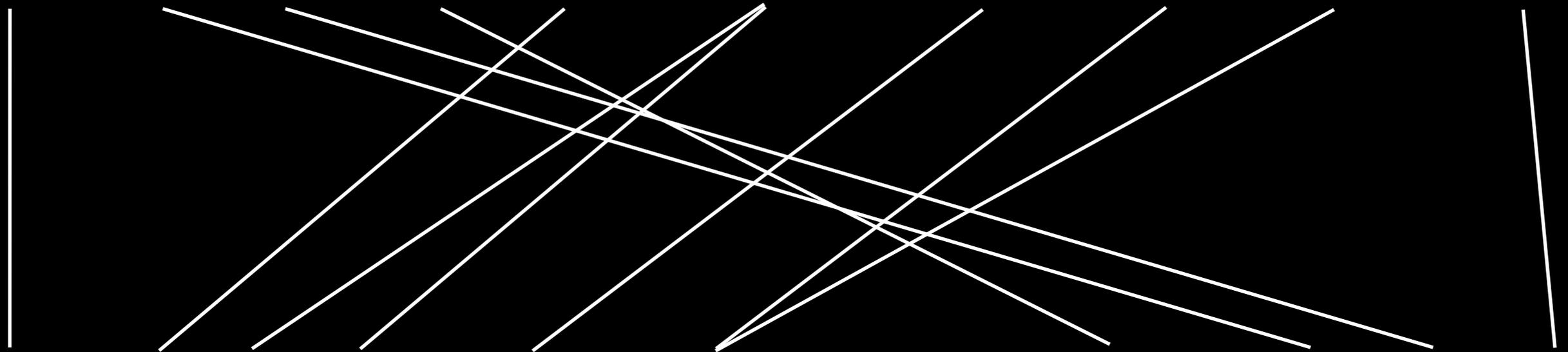
However , the sky remained clear under the strong north wind .

$$p(\textit{however} | \text{虽然}) = \frac{\text{\# of times 虽然 aligns to However}}{\text{\# of times 虽然 occurs}}$$

IBM Model 1

Although north wind howls , but sky still very clear .

虽然 北 风 呼啸 ， 但 天空 依然 十分 清澈 。



However , the sky remained clear under the strong north wind .

$$p(\textit{however} | \text{虽然}) = \frac{\# \text{ of times 虽然 aligns to However}}{\# \text{ of times 虽然 occurs}}$$

A man in a dark suit, white shirt, and red tie, wearing glasses, is pointing his right index finger towards the right. He is positioned on the left side of the frame. The background is a blue graphic with a world map and a row of stars at the bottom. The text 'THE WORD' is prominently displayed in the upper right quadrant.

THE WORD

- **Optimization**

MLE for IBM Model 1 (observed)

$$\hat{\theta} = \arg \max_{\theta} \prod_{n=1}^N \left(p(I^{(n)} | J^{(n)}) \prod_{i=1}^{I^{(n)}} p(a_i^{(n)} | J^{(n)}) \cdot p(f_i^{(n)} | e_{a_i^{(n)}}) \right)$$

MLE for IBM Model 1 (observed)

number of
sentences

alignment of French
word at position i

$$\hat{\theta} = \arg \max_{\theta} \prod_{n=1}^N \left(p(I^{(n)} | J^{(n)}) \prod_{i=1}^{I^{(n)}} p(a_i^{(n)} | J^{(n)}) \cdot p(f_i^{(n)} | e_{a_i^{(n)}}) \right)$$

French, English
sentence lengths

French, English
word pair

MLE for IBM Model 1 (observed)

$$\hat{\theta} = \arg \max_{\theta} \prod_{n=1}^N \left(\underbrace{p(I^{(n)} | J^{(n)}) \prod_{i=1}^{I^{(n)}} p(a_i^{(n)} | J^{(n)})}_{\text{constant!}} \cdot p(f_i^{(n)} | e_{a_i}^{(n)}) \right)$$

MLE for IBM Model 1 (observed)

$$\hat{\theta} = \arg \max_{\theta} C \prod_{n=1}^N \prod_{i=1}^{I^{(n)}} p(f_i^{(n)} | e_{a_i}^{(n)})$$

MLE for IBM Model 1 (observed)

$$\hat{\theta} = \arg \max_{\theta} \log \left(C \prod_{n=1}^N \prod_{i=1}^{I^{(n)}} p(f_i^{(n)} | e_{a_i}^{(n)}) \right)$$

$$\log(a) < \log(b) \iff a < b$$

MLE for IBM Model 1 (observed)

$$\hat{\theta} = \arg \max_{\theta} \log \left(C \cdot \prod_{f,e} p(f|e)^{\text{count}(\langle f,e \rangle)} \right)$$

MLE for IBM Model 1 (observed)

$$\hat{\theta} = \arg \max_{\theta} \log C + \sum_{f,e} \text{count}(\langle f, e \rangle) \log p(f|e)$$

log of product = sum of logs

MLE for IBM Model 1 (observed)

$$\Lambda(\theta, \lambda) = \log C + \sum_{f,e} \text{count}(\langle f, e \rangle) \log p(f|e) - \sum_e \lambda_e \underbrace{\left(\sum_f p(f|e) - 1 \right)}$$

Lagrange multiplier expresses normalization constraint

MLE for IBM Model 1 (observed)

$$\Lambda(\theta, \lambda) = \log C + \sum_{f,e} \text{count}(\langle f, e \rangle) \log p(f|e) - \sum_e \lambda_e \left(\sum_f p(f|e) - 1 \right)$$

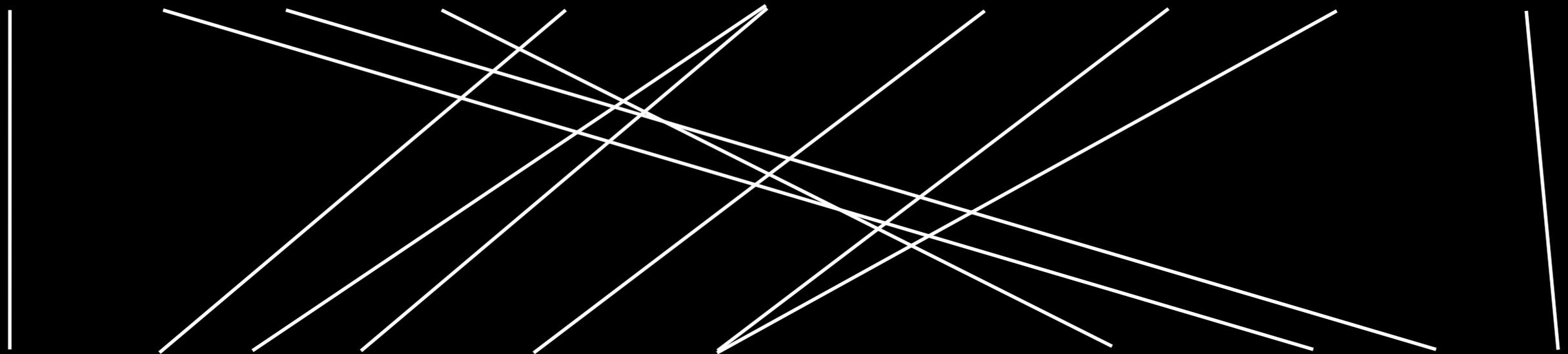
derivative

$$\frac{\partial \Lambda(\theta, \lambda)}{\partial p(f|e)} = \frac{\text{count}(\langle f, e \rangle)}{p(f|e)} - \lambda_e$$

MLE for IBM Model 1 (observed)

Although north wind howls , but sky still very clear .

虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。



However , the sky remained clear under the strong north wind .

$$p(\textit{however} | \text{虽然}) = \frac{\# \text{ of times 虽然 aligns to However}}{\# \text{ of times 虽然 occurs}}$$

MLE for IBM Model 1 (unobserved)

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。

However , the sky remained clear under the strong north wind .

$$p(\textit{however} | \text{虽然}) = ???$$

MLE for IBM Model 1 (observed)

$$\hat{\theta} = \arg \max_{\theta} \log \left(C \prod_{n=1}^N \prod_{i=1}^{I^{(n)}} p(f_i^{(n)} | e_{a_i}^{(n)}) \right)$$

MLE for IBM Model 1 (unobserved)

$$\hat{\theta} = \arg \max_{\theta} \log \left(C \prod_{n=1}^N \sum_a \prod_{i=1}^{I^{(n)}} p(f_i^{(n)} | e_{a_i}^{(n)}) \right)$$

marginalize over alignments:

$$p(f|e) = \sum_a p(f, a|e)$$

MLE for IBM Model 1 (unobserved)

$$\hat{\theta} = \arg \max_{\theta} \log \left(C \cdot \prod_{f,e} p(f|e)^{\mathbb{E}[\text{count}(\langle f,e \rangle)]} \right)$$

MLE for IBM Model 1 (unobserved)

$$\hat{\theta} = \arg \max_{\theta} \log \left(C \cdot \prod_{f,e} p(f|e)^{\mathbb{E}[\text{count}(\langle f,e \rangle)]} \right)$$

Not constant! Depends on parameters,
no analytic solution.



MLE for IBM Model 1 (unobserved)

$$\hat{\theta} = \arg \max_{\theta} \log \left(C \cdot \prod_{f,e} p(f|e)^{\mathbb{E}[\text{count}(\langle f,e \rangle)]} \right)$$

Not constant! Depends on parameters,
no analytic solution.

But it does strongly imply an iterative solution.

Likelihood Estimation for Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

Parameters and alignments are both unknown.

However , the sky remained clear under the strong north wind .

$p(\textit{English word}|\textit{Chinese word})$ unobserved!

Likelihood Estimation for Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

Parameters and alignments are both unknown.

If we knew the alignments, we could calculate the values of the parameters.

However , the sky remained clear under the strong north wind .

$p(\text{English word}|\text{Chinese word})$ unobserved!

Likelihood Estimation for Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

Parameters and alignments are both unknown.

If we knew the alignments, we could calculate the values of the parameters.

If we knew the parameters, we could calculate the likelihood of the data.

However , the sky remained clear under the strong north wind .

$p(\text{English word}|\text{Chinese word})$ unobserved!

Likelihood Estimation for Model 1

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

Parameters and alignments are both unknown.

If we knew the alignments, we could calculate the values of the parameters.



If we knew the parameters, we could calculate the likelihood of the data.

However , the sky remained clear under the strong north wind .

$p(\text{English word}|\text{Chinese word})$ unobserved!

The Plan: Bootstrapping

- Arbitrarily select a set of parameters (say, uniform).
- Calculate *expected counts* of the unseen events.
- Choose new parameters to maximize likelihood, using expected counts as proxy for observed counts.
- Iterate.
- Guarantee: likelihood will be monotonically nondecreasing.

The Plan: Bootstrapping

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

However , the sky remained clear under the strong north wind .

The Plan: Bootstrapping

Although north wind howls , but sky still very clear .

虽然北风呼啸，但天空依然十分清澈。ε

if we had observed the alignment, this line would either be here (count 1) or it wouldn't (count 0).

However , the sky remained clear under the strong north wind .

The Plan: Bootstrapping

Although north wind howls , but sky still very clear .

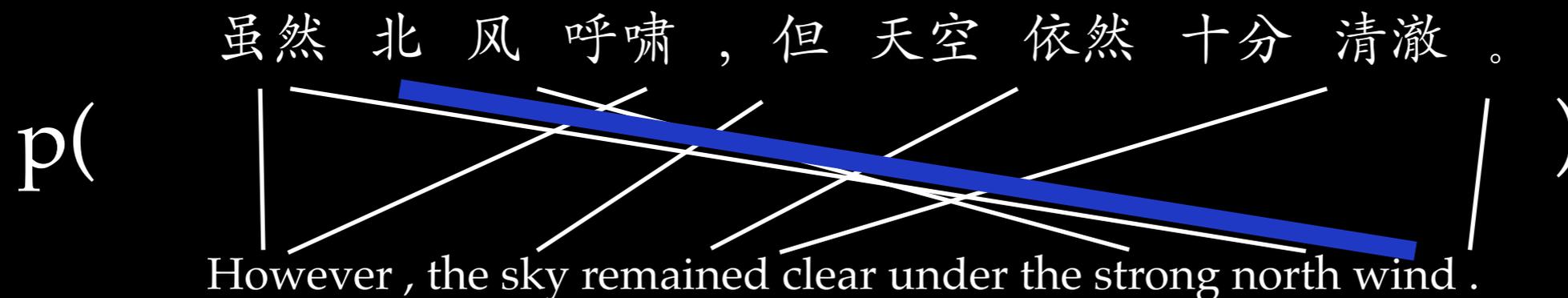
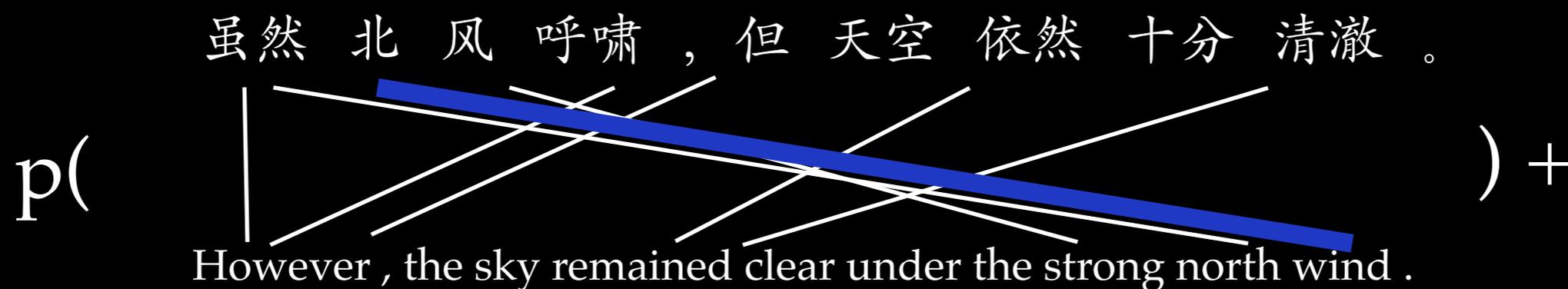
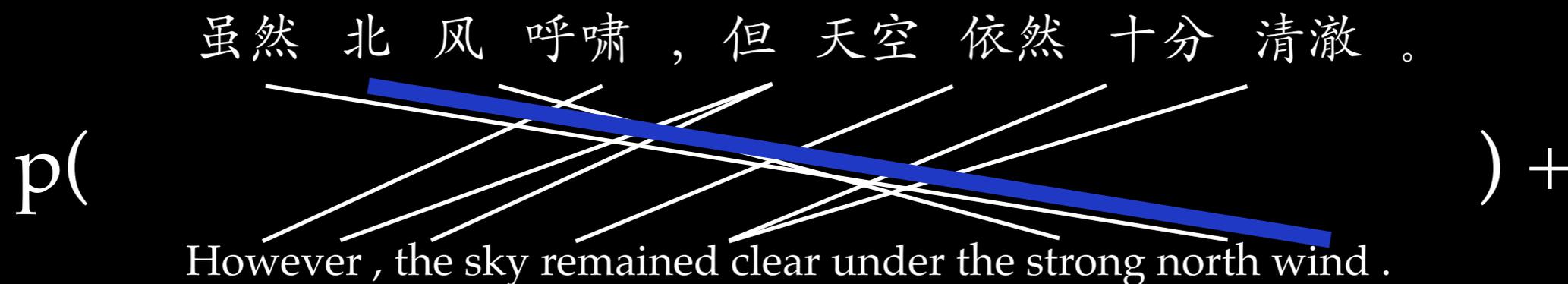
虽然北风呼啸，但天空依然十分清澈。 ϵ

if we had observed the alignment, this line would either be here (count 1) or it wouldn't (count 0).

since we didn't observe the alignment, we calculate the probability that it's there.

However , the sky remained clear under the strong north wind .

Marginalize: sum all alignments containing the link



Divide by sum of all *possible* alignments

虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。

p(

However , the sky remained clear under the strong north wind .

) +

虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。

p(

However , the sky remained clear under the strong north wind .

) +

虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。

p(

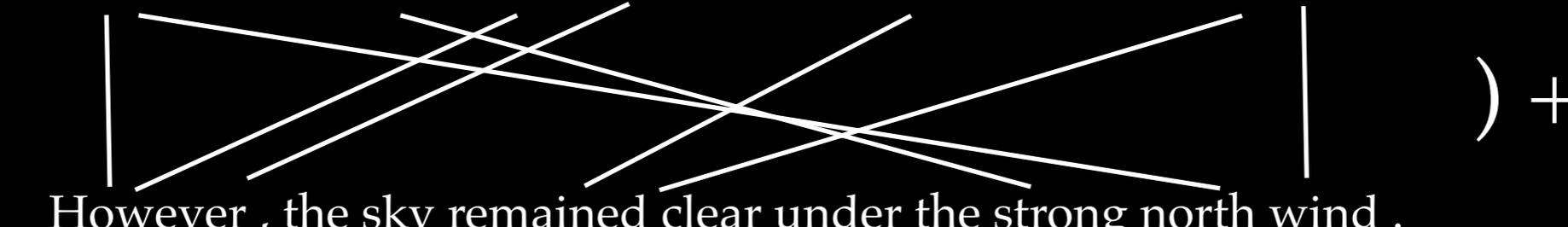
However , the sky remained clear under the strong north wind .

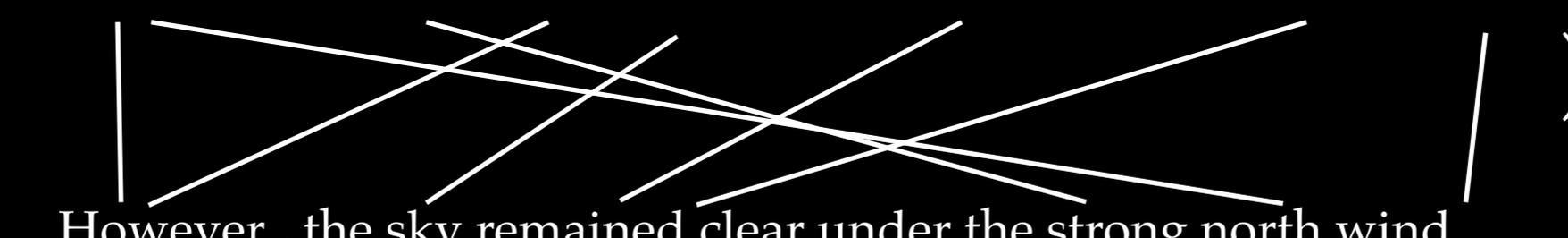
)

Divide by sum of all *possible* alignments

$p(\text{虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。}) +$

 $\text{However , the sky remained clear under the strong north wind .}$

$p(\text{虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。}) +$

 $\text{However , the sky remained clear under the strong north wind .}$

$p(\text{虽然 北 风 呼 啸 ， 但 天 空 依 然 十 分 清 澈 。})$

 $\text{However , the sky remained clear under the strong north wind .}$

Is this hard? How many alignments are there?

Expectation Maximization

probability of an alignment.

$$p(F, A|E) = p(I|J) \prod_{a_i} p(a_i = j) p(f_i|e_j)$$

Expectation Maximization

probability of an alignment.

$$p(F, A|E) = p(I|J) \prod_{a_i} p(a_i = j) p(f_i|e_j)$$

observed uniform

Expectation Maximization

probability of an alignment.

factors across words.

$$p(F, A|E) = p(I|J) \prod_{a_i} p(a_i = j) p(f_i|e_j)$$

observed

uniform

Expectation Maximization

marginal probability of
alignments containing link

$$\sum_{a \in A: \text{北} \leftrightarrow \text{north}} p(\text{north} | \text{北}) \cdot p(\text{rest of } a)$$

Expectation Maximization

marginal probability of
alignments containing link

$$p(\textit{north} | \text{北}) = \sum_{a \in A: \text{北} \leftrightarrow \textit{north}} p(\textit{rest of } a)$$

Expectation Maximization

marginal probability of
alignments containing link

$$p(\textit{north} | \text{北}) \sum_{a \in A: \text{北} \leftrightarrow \textit{north}} p(\textit{rest of } a)$$

$$\sum_{c \in \textit{Chinese words}} p(\textit{north} | c) \sum_{a \in A: c \leftrightarrow \textit{north}} p(\textit{rest of } a)$$

marginal probability of all
alignments

Expectation Maximization

marginal probability of
alignments containing link

$$p(\textit{north} | \text{北}) \sum_{a \in A: \text{北} \leftrightarrow \textit{north}} p(\textit{rest of } a)$$

$$\sum_{c \in \textit{Chinese words}} p(\textit{north} | c) \sum_{a \in A: c \leftrightarrow \textit{north}} p(\textit{rest of } a)$$

identical!

marginal probability of all
alignments

Expectation Maximization

$$\frac{p(\textit{north} | \text{北})}{\sum_{c \in \textit{Chinese words}} p(\textit{north} | c)}$$

Expectation Maximization

marginal probability (expected count) of an alignment containing the link

$$\frac{p(\textit{north} | \text{北})}{\sum_{c \in \textit{Chinese words}} p(\textit{north} | c)}$$

Expectation Maximization

marginal probability (expected count) of an alignment containing the link

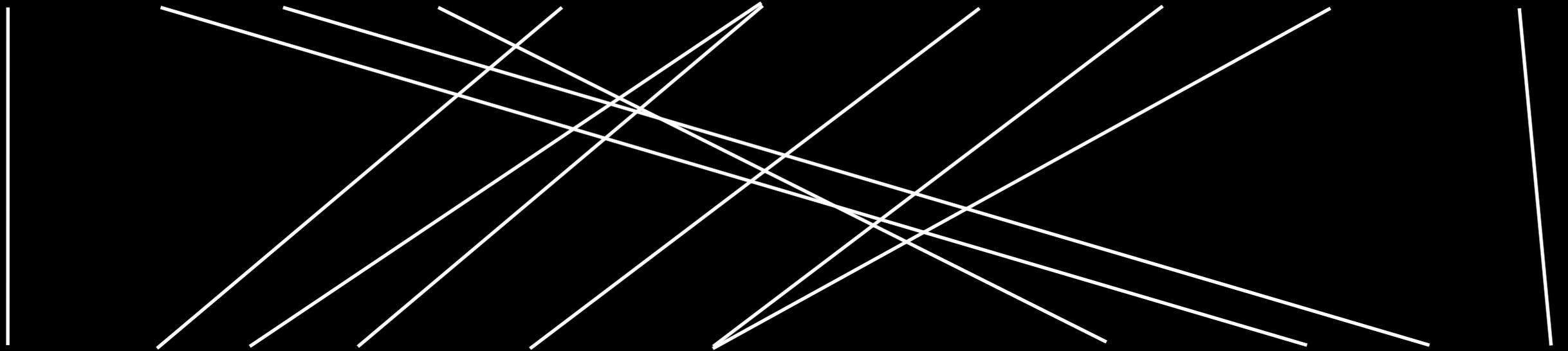
$$\frac{p(\textit{north} | \text{北})}{\sum_{c \in \textit{Chinese words}} p(\textit{north} | c)}$$

For each sentence, use this quantity instead of 0 or 1

Translation Models

Although north wind howls , but sky still very clear .

虽然 北 风 呼啸 ， 但 天空 依然 十分 清澈 。



However , the sky remained clear under the strong north wind .

$$p(\textit{however} | \text{虽然}) = \frac{\# \text{ of times 虽然 aligns to However}}{\# \text{ of times 虽然 occurs}}$$

Translation Models

Although north wind howls , but sky still very clear .

虽然 北 风 呼啸 ， 但 天空 依然 十分 清澈 。

However , the sky remained clear under the strong north wind .

$$p(\textit{however} | \text{虽然}) = \frac{\textit{Expected} \# \text{ of times 虽然 aligns to However}}{\# \text{ of times 虽然 occurs}}$$

Expectation Maximization

Why does this even work?

$$\frac{p(\textit{north} | \text{北})}{\sum_{c \in \textit{Chinese words}} p(\textit{north} | c)}$$

Expectation Maximization

Observation 1: We are still solving a maximum likelihood estimation problem.

Expectation Maximization

Observation 1: We are still solving a maximum likelihood estimation problem.

$$p(\textit{Chinese}|\textit{English}) = \sum_{\textit{alignments}} p(\textit{Chinese}, \textit{alignment}|\textit{English})$$

Expectation Maximization

Observation 1: We are still solving a maximum likelihood estimation problem.

$$p(\textit{Chinese}|\textit{English}) = \sum_{\textit{alignments}} p(\textit{Chinese}, \textit{alignment}|\textit{English})$$

MLE: choose parameters that maximize this expression.

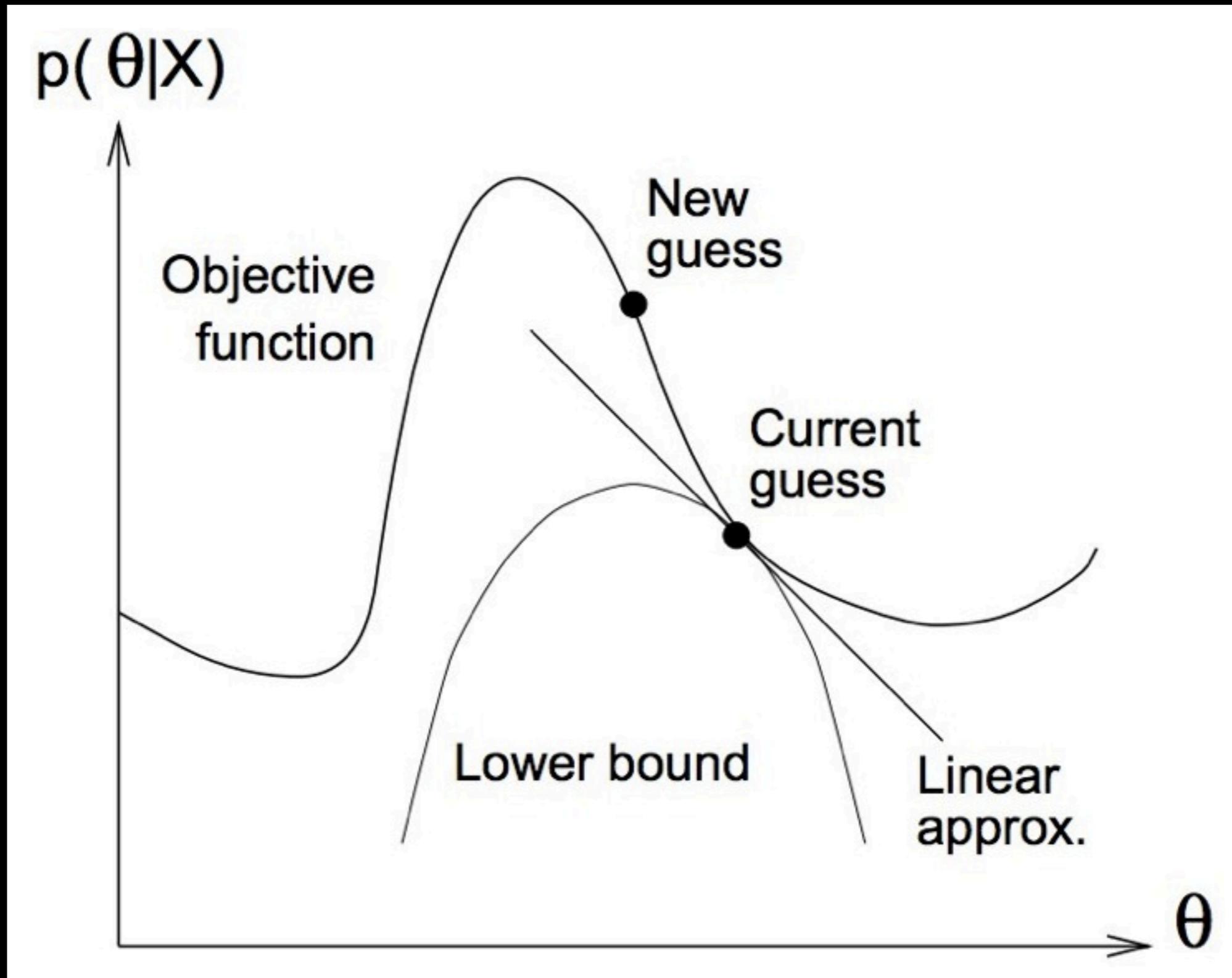
Expectation Maximization

Observation 1: We are still solving a maximum likelihood estimation problem.

$$p(\textit{Chinese}|\textit{English}) = \sum_{\textit{alignments}} p(\textit{Chinese}, \textit{alignment}|\textit{English})$$

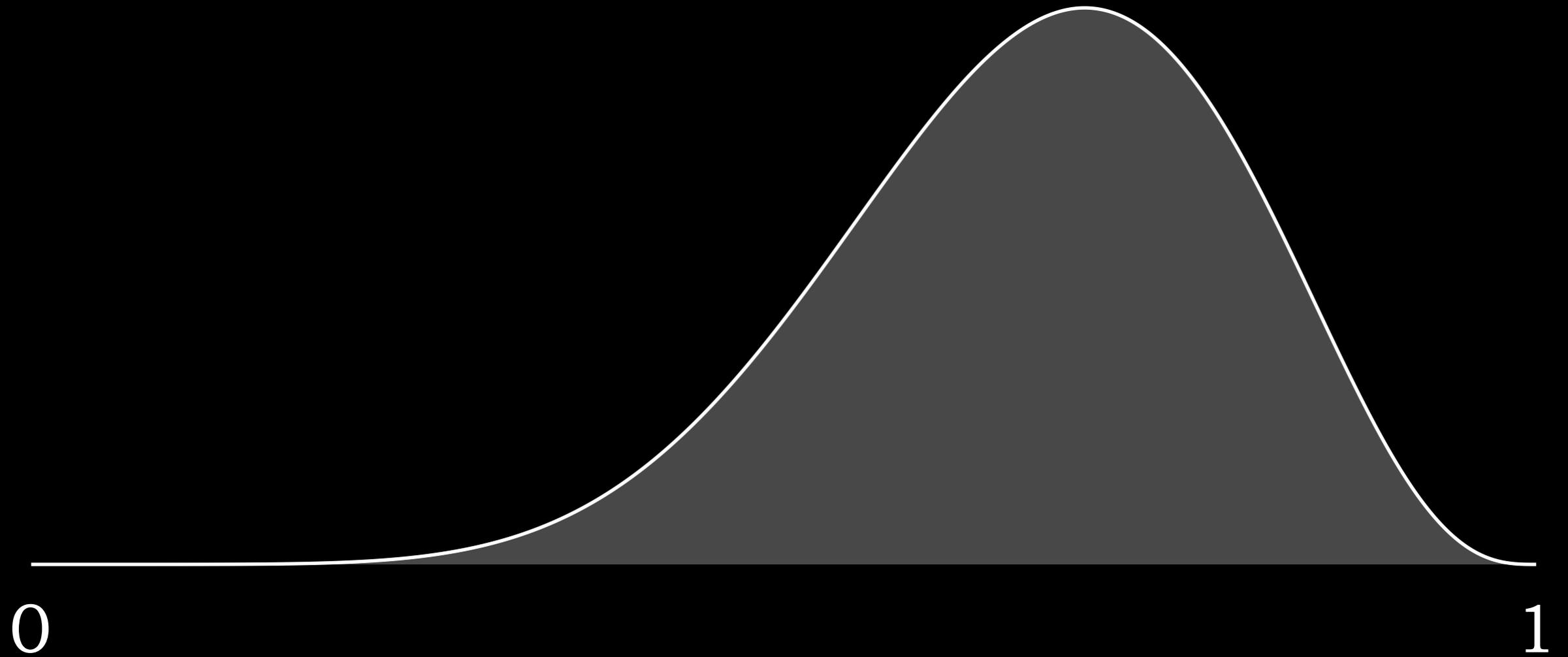
MLE: choose parameters that maximize this expression.

Minor problem: there is no analytic solution.



(from Minka '98)

... and, likelihood is *convex* for this model:



Summary

- Many possible models.
- Many possible objective functions.
- *Learning is optimization*: choose parameters that optimize some function, such as likelihood.
- Try some out this afternoon in the lab!