

# Machine Translation

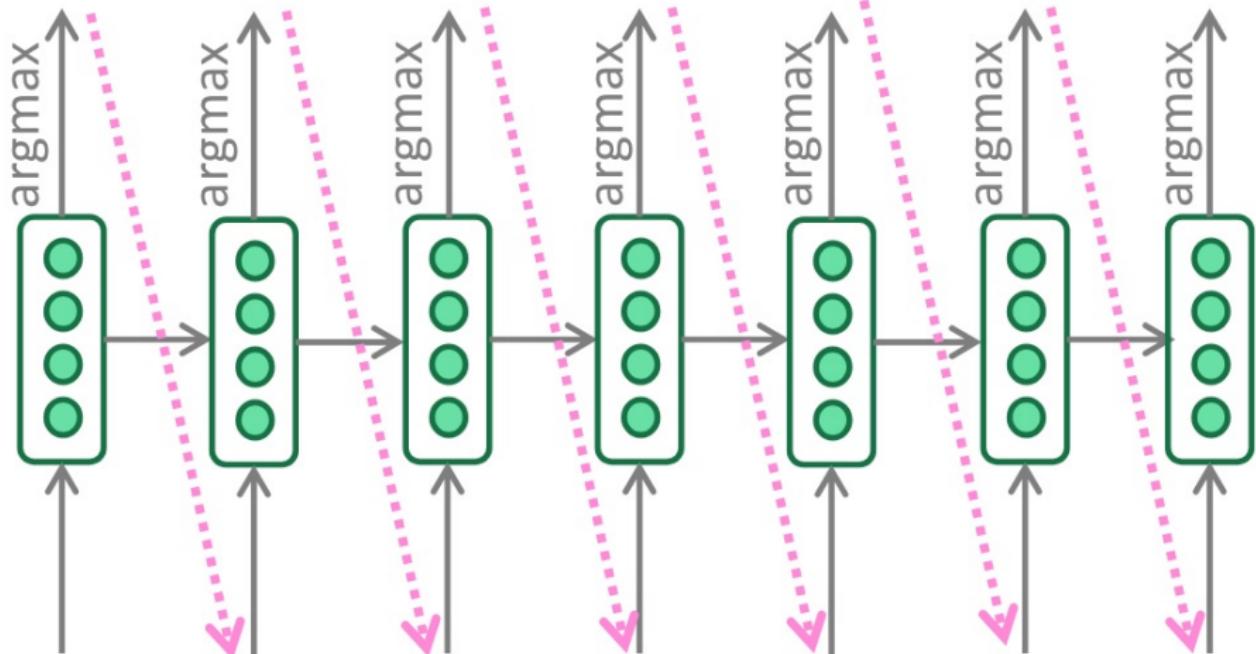
Jordan Boyd-Graber

University of Maryland

Decoding

Adapted from material by Mohit Iyyer, Luke Zettlemoyer, Kalpesh Krishna, Karthik Narasimhan, Greg Durrett, Chris Manning, Dan Jurafsky

the poor don't have any money <END>



Argmax at every time step

# Sampling Methods

Softmax distribution

$$p(w) = \frac{\exp\{\beta \cdot \vec{f}(w)\}}{\sum_{w'} \exp\{\beta \cdot \vec{f}(w')\}} \quad (1)$$

- top- $k$
- Nucleus / top- $p$
- Temperature

# Sampling Methods

## Softmax distribution

$$p(w) = \frac{\exp\{\beta \cdot \vec{f}(w)\}}{\sum_{w'} \exp\{\beta \cdot \vec{f}(w')\}} \quad (1)$$

- top- $k$ : Only sample from  $k$  items with highest probability
- Nucleus / top- $p$
- Temperature

# Sampling Methods

Softmax distribution

$$p(w) = \frac{\exp\{\beta \cdot \vec{f}(w)\}}{\sum_{w'} \exp\{\beta \cdot \vec{f}(w')\}} \quad (1)$$

- top- $k$ : Only sample from  $k$  items with highest probability
- Nucleus / top- $p$ : Only sample from highest items with at least  $p$  probability
- Temperature

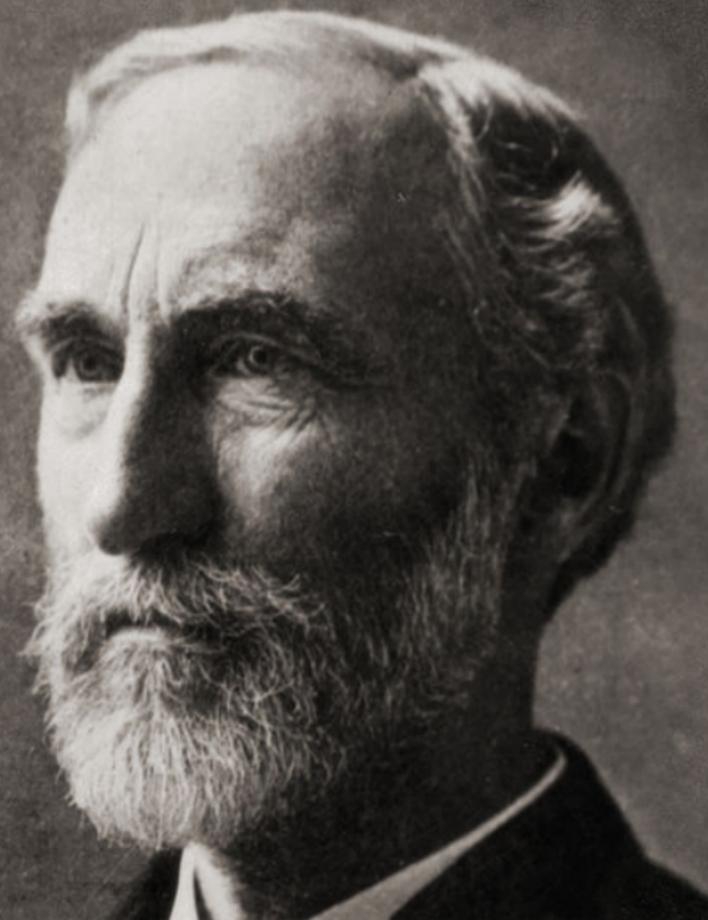
# Sampling Methods

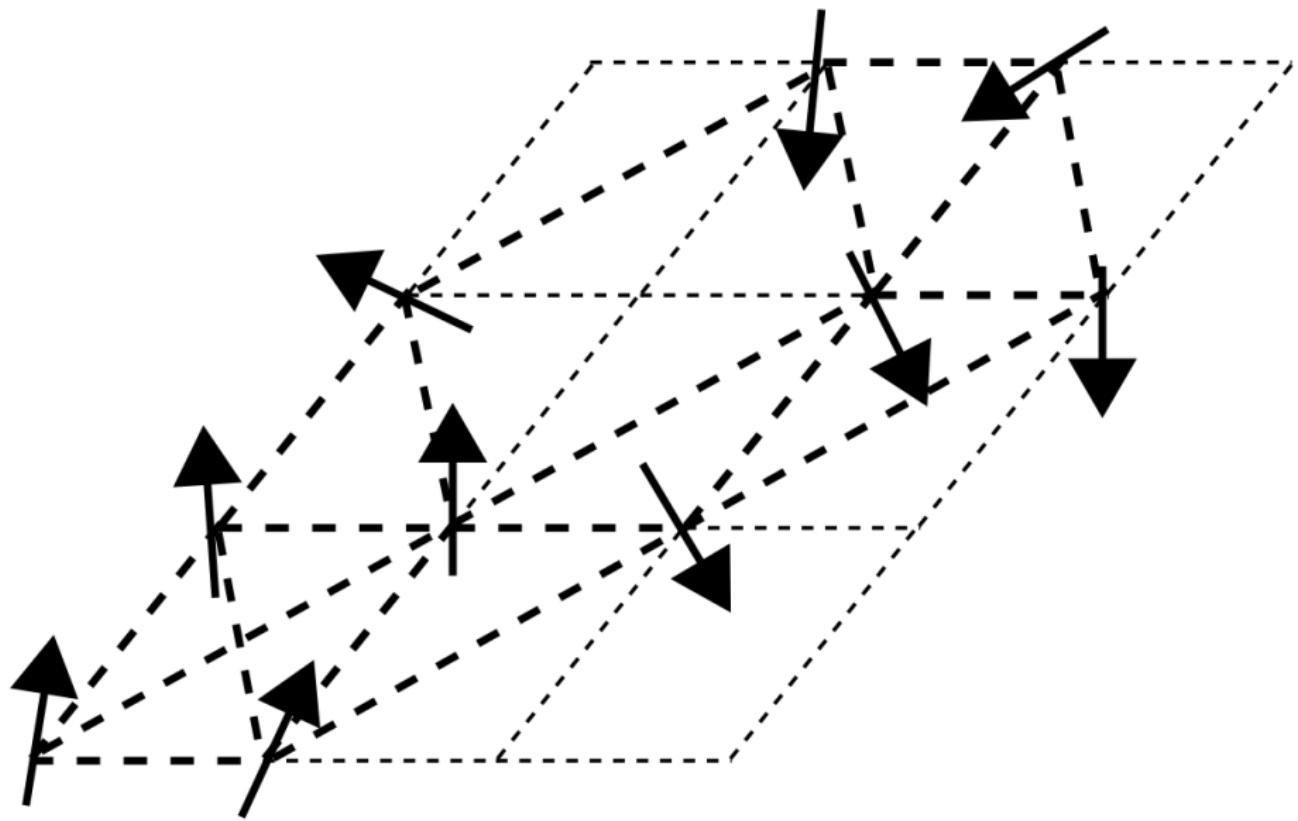
Softmax distribution

$$p(w) = \frac{\exp\{\beta \cdot \vec{f}(w)\}}{\sum_{w'} \exp\{\beta \cdot \vec{f}(w')\}} \quad (1)$$

- top- $k$ : Only sample from  $k$  items with highest probability
- Nucleus / top- $p$ : Only sample from highest items with at least  $p$  probability
- Temperature

$$p(w) = \frac{\exp\left\{\frac{\beta \cdot \vec{f}(w)}{T}\right\}}{\sum_{w'} \exp\left\{\frac{\beta \cdot \vec{f}(w')}{T}\right\}} \quad (2)$$





# Sampling Methods

Softmax distribution

$$p(w) = \frac{\exp\{\beta \cdot \vec{f}(w)\}}{\sum_{w'} \exp\{\beta \cdot \vec{f}(w')\}} \quad (3)$$

- top- $k$
- Nucleus / top- $p$
- Temperature

# Sampling Methods

## Softmax distribution

$$p(w) = \frac{\exp\{\beta \cdot \vec{f}(w)\}}{\sum_{w'} \exp\{\beta \cdot \vec{f}(w')\}} \quad (3)$$

- top- $k$ : Only sample from  $k$  items with highest probability
- Nucleus / top- $p$
- Temperature

## Sampling Methods

Softmax distribution

$$p(w) = \frac{\exp\{\beta \cdot \vec{f}(w)\}}{\sum_{w'} \exp\{\beta \cdot \vec{f}(w')\}} \quad (3)$$

- top- $k$ : Only sample from  $k$  items with highest probability
- Nucleus / top- $p$ : Only sample from highest items with at least  $p$  probability
- Temperature

# Sampling Methods

Softmax distribution

$$p(w) = \frac{\exp\{\beta \cdot \vec{f}(w)\}}{\sum_{w'} \exp\{\beta \cdot \vec{f}(w')\}} \quad (3)$$

- top- $k$ : Only sample from  $k$  items with highest probability
- Nucleus / top- $p$ : Only sample from highest items with at least  $p$  probability
- Temperature

$$p(w) = \frac{\exp\left\{\frac{\beta \cdot \vec{f}(w)}{T}\right\}}{\sum_{w'} \exp\left\{\frac{\beta \cdot \vec{f}(w')}{T}\right\}} \quad (4)$$

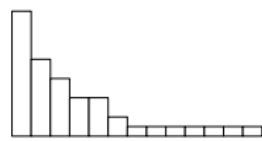
Top- $k$

Nucleus

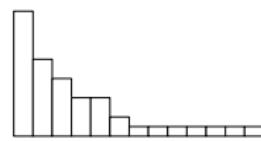
Temperature (T=0.1)

Temperature (T=2)

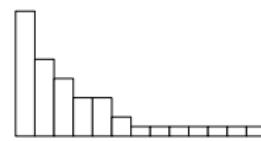
Top- $k$



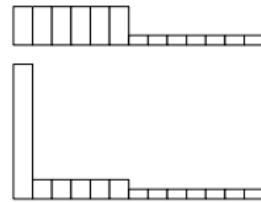
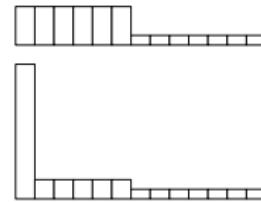
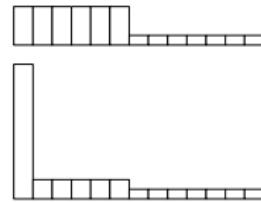
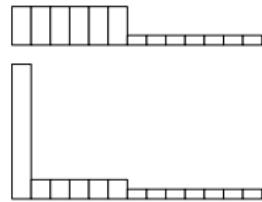
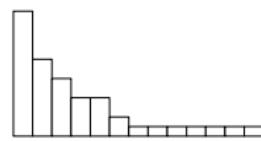
Nucleus



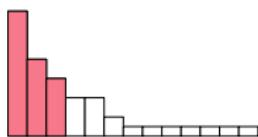
Temperature (T=0.1)



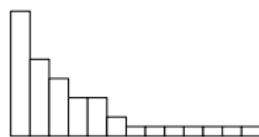
Temperature (T=2)



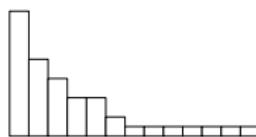
Top- $k$



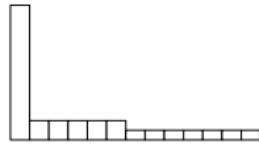
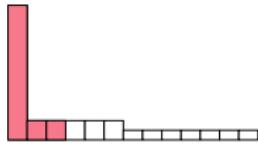
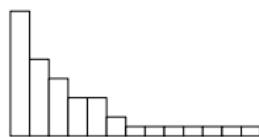
Nucleus



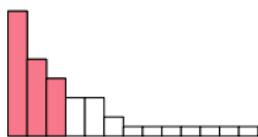
Temperature (T=0.1)



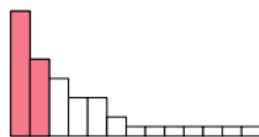
Temperature (T=2)



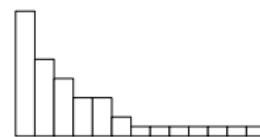
Top- $k$



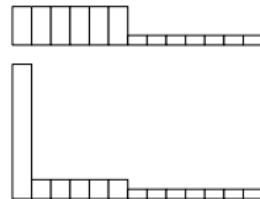
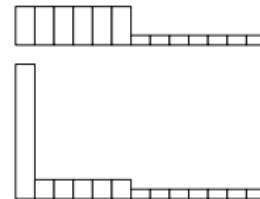
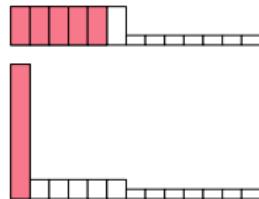
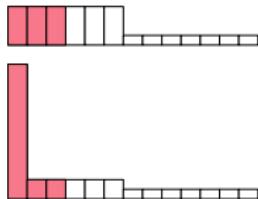
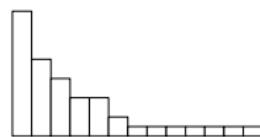
Nucleus

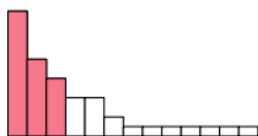


Temperature (T=0.1)

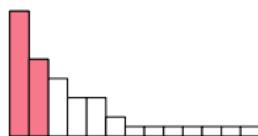


Temperature (T=2)

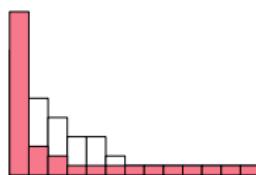


Top- $k$ 

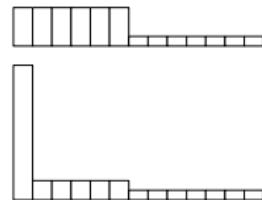
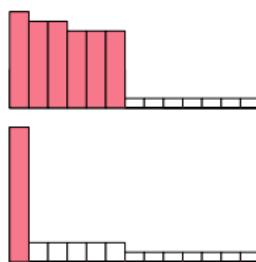
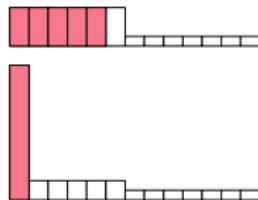
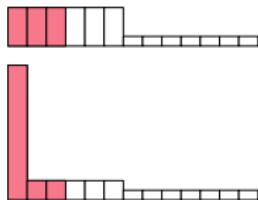
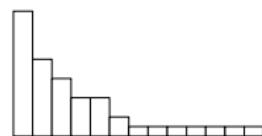
Nucleus

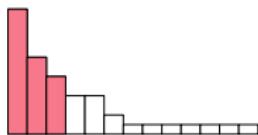


Temperature (T=0.1)

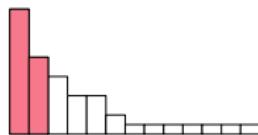


Temperature (T=2)

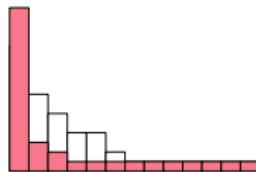


Top- $k$ 

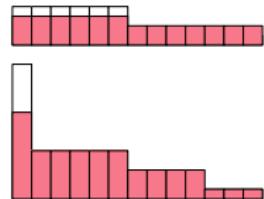
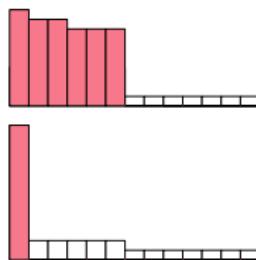
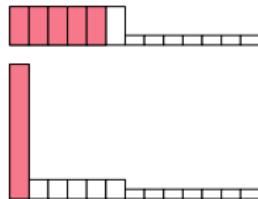
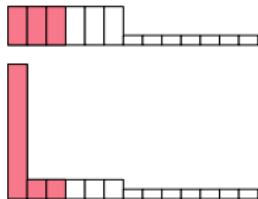
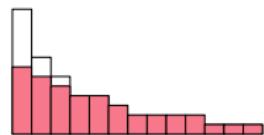
Nucleus



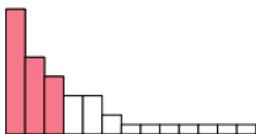
Temperature (T=0.1)



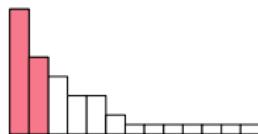
Temperature (T=2)



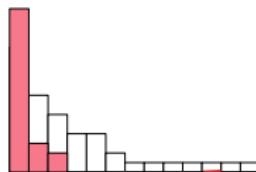
Top- $k$



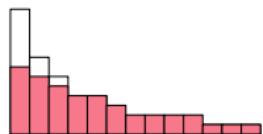
Nucleus



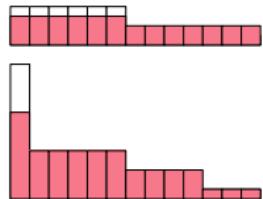
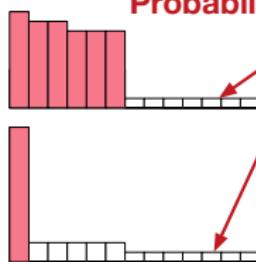
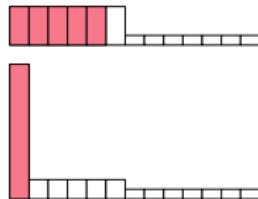
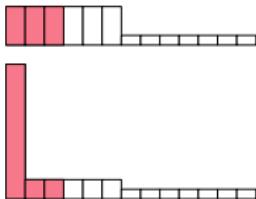
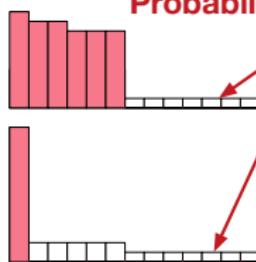
Temperature ( $T=0.1$ )



Temperature ( $T=2$ )



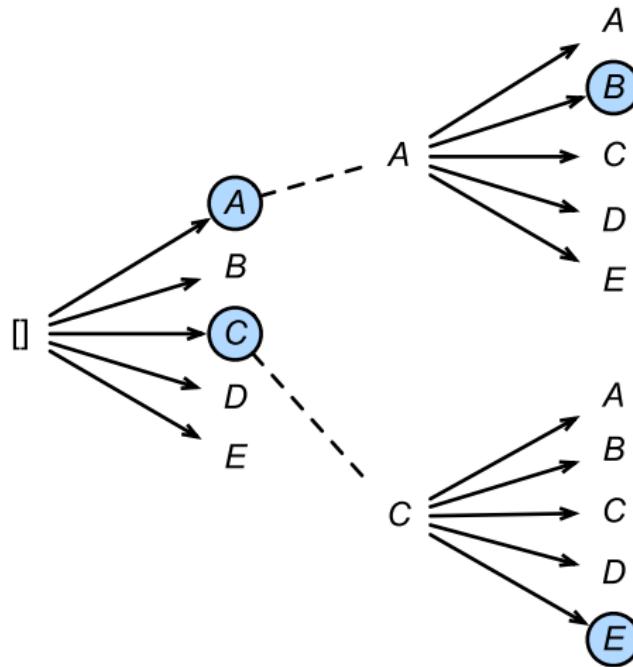
Probabilities too small to see!



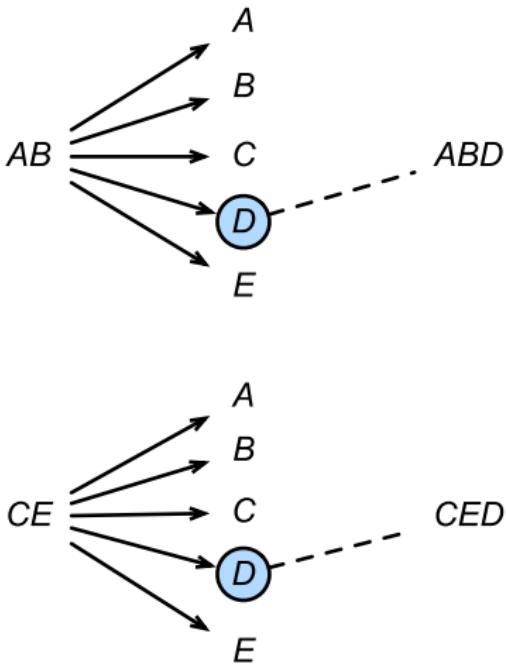
## What do you do with samples?

- Getting out of being stuck in a garden path
- Getting diverse outputs
- Combining multiple models together
- Rescoring by a non-probability metric

Time step 1  
Candidates



Time step 2  
Candidates



Time step 3  
Candidates

From Zhang et al. (Dive into Deep Learning)

<S> —→

—→ </S>

—→ </S>

—→ </S>

—→ </S>

<S> →

→ </S>



→ </S>

→ </S>

<S> →



→ </S>

→ </S>



→ </S>

→ </S>

< s > →



→ < /s >

→ < /s >

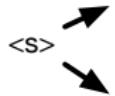


→ < /s >

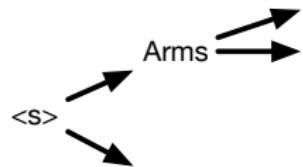
→ < /s >

<s>

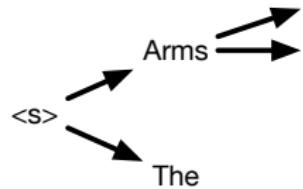
Beam Search Decoding for: Die Arme haben kein Geld



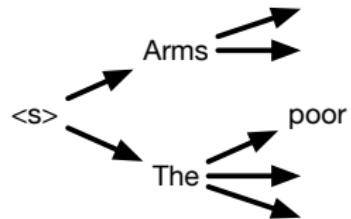
Beam Search Decoding for: Die Arme haben kein Geld



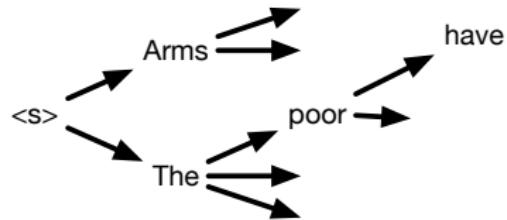
Beam Search Decoding for: Die Arme haben kein Geld



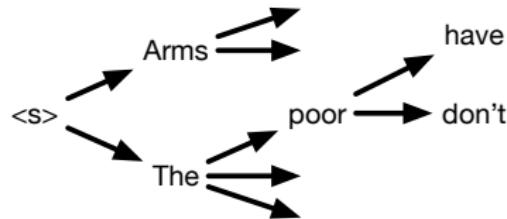
Beam Search Decoding for: Die Arme haben kein Geld



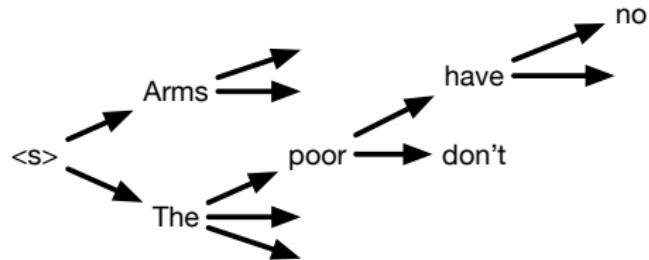
Beam Search Decoding for: Die Arme haben kein Geld



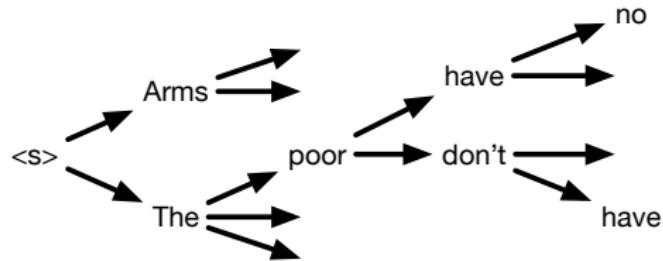
Beam Search Decoding for: Die Arme haben kein Geld



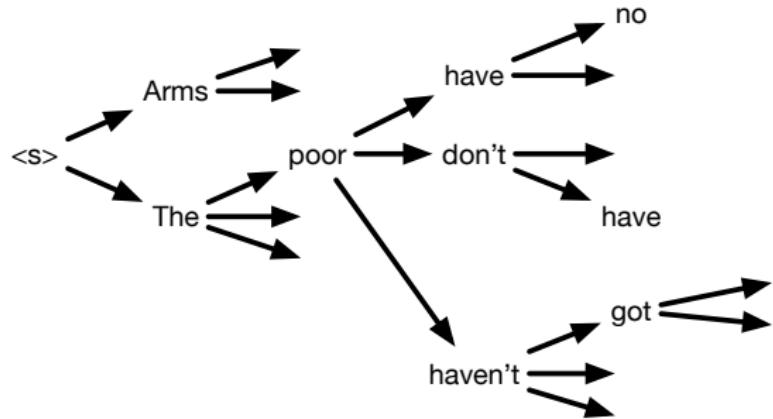
Beam Search Decoding for: Die Arme haben kein Geld



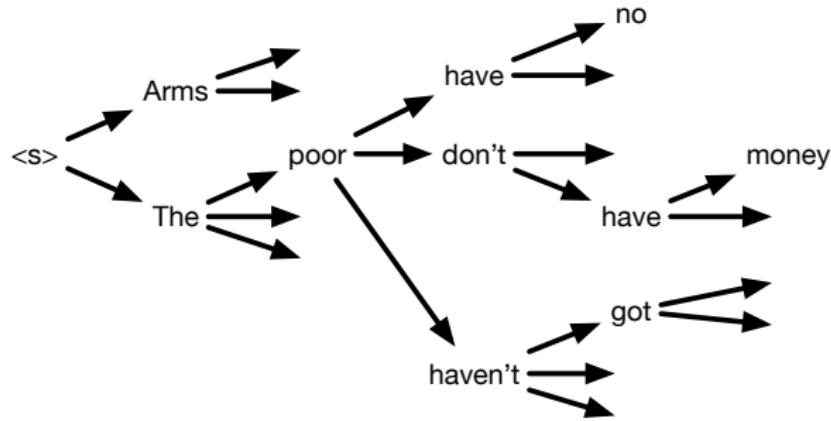
Beam Search Decoding for: Die Arme haben kein Geld



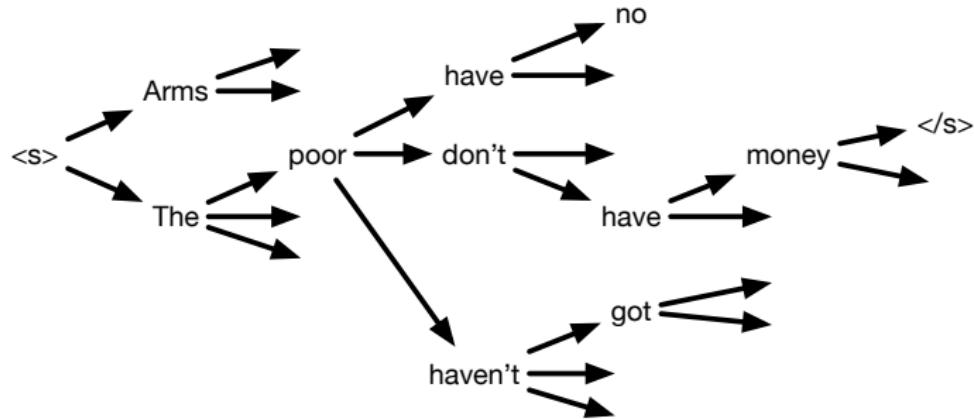
Beam Search Decoding for: Die Arme haben kein Geld



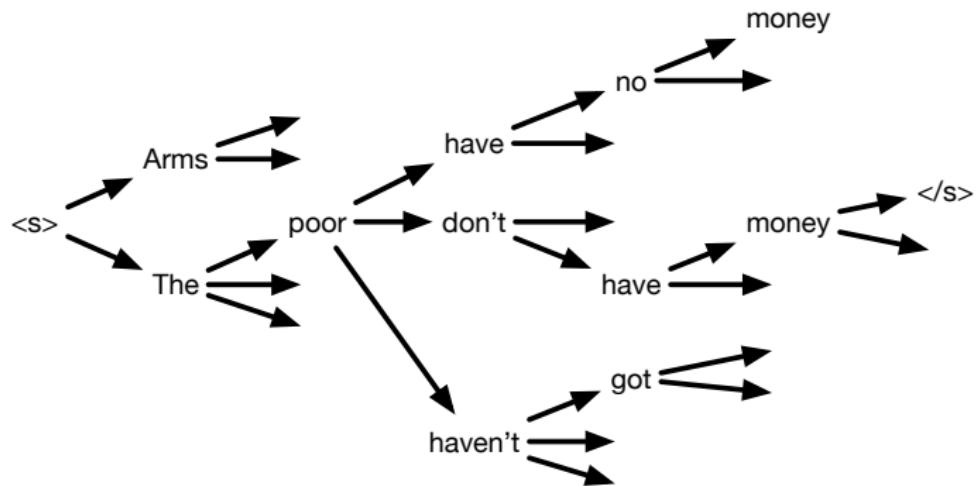
Beam Search Decoding for: Die Arme haben kein Geld



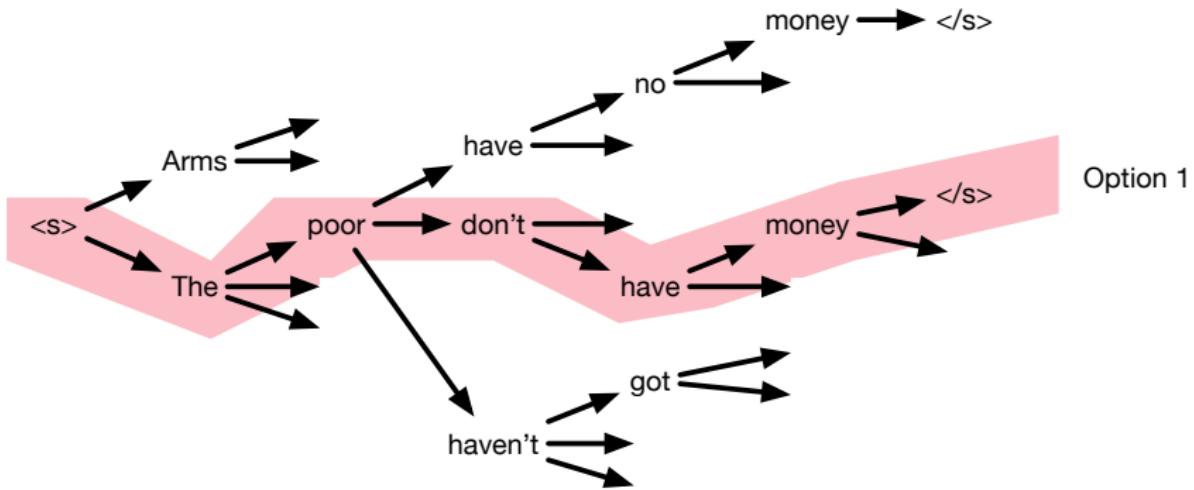
Beam Search Decoding for: Die Arme haben kein Geld



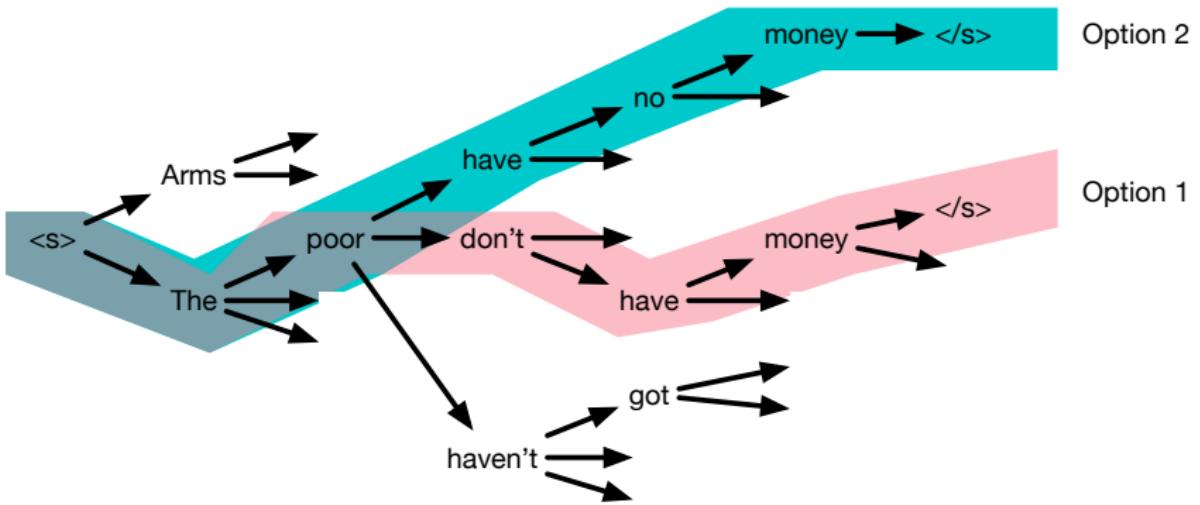
Beam Search Decoding for: Die Arme haben kein Geld



Beam Search Decoding for: Die Arme haben kein Geld



Beam Search Decoding for: Die Arme haben kein Geld



Beam Search Decoding for: Die Arme haben kein Geld

## Using multiple sources

- Generate from multiple models
- Generate from multiple directions
- Generate from multiple data
- Generate from multiple temperatures

## How to pick?

- Show to a user

## How to pick?

- Show to a user
- Take highest probability

# How to pick?

## Reinforcement Learning for Bandit Neural Machine Translation with Simulated Human Feedback

Khanh Nguyen<sup>○◊</sup> and Hal Daumé III<sup>○◆◊▽</sup> and Jordan Boyd-Graber<sup>○▲◆◊</sup>

University of Maryland: Computer Science,<sup>○</sup>Language Science,<sup>◆</sup>iSchool,<sup>▲</sup>UMIACS<sup>◊</sup>

Microsoft Research, New York<sup>○</sup>

{kxnguyen, hal, jbg}@umiacs.umd.edu

## Can Neural Machine Translation be Improved with User Feedback?

Julia Kreutzer<sup>1,\*</sup> and Shahram Khadivi<sup>3</sup> and Evgeny Matusov<sup>2</sup> and Stefan Riezler<sup>1,2</sup>

<sup>1</sup>Computational Linguistics & <sup>2</sup>IWR, Heidelberg University, Germany

{kreutzer, riezler}@cl.uni-heidelberg.de

<sup>3</sup>eBay Inc., Aachen, Germany

{skhadivi, ematusov}@eBay.com

- Show to a user
- Take highest probability
- Rerank



## RankGen — Improving Text Generation with Large Ranking Models (*EMNLP 2022*)



Kalpesh Krishna



Yapei Chang



John Wieting



Mohit Iyyer

UMassAmherst  
Manning College of Information  
& Computer Sciences



Decoding Method	GPT2-md		GPT2-XL	
	PG19	wiki	PG19	wiki
Nucleus ( $p = 0.9$ )	73.0	74.6	74.4	75.0
Eta ( <a href="#">Hewitt et al., 2022</a> )	76.9	71.2	76.9	74.8
<i>Contrastive methods</i>				
search ( <a href="#">Su et al., 2022</a> )	5.3	21.2	54.0	43.2
decode ( <a href="#">Li et al., 2022</a> )	65.2	83.2	73.2	84.9
RANKGEN-all-XL (ours)				
rerank full ancestral	<b>79.0</b>	84.9	<b>79.0</b>	86.4
beam search nucleus	76.2	<b>88.9</b>	77.0	<b>89.4</b>

# Automatic Song Translation for Tonal Languages

**Fenfei Guo**

University of Maryland

[fguo1@umd.edu](mailto:fguo1@umd.edu)

**Qixin He**

Purdue University

[heqixin@purdue.edu](mailto:heqixin@purdue.edu)

**Chen Zhang**

Zhejiang University

[zc99@zju.edu.cn](mailto:zc99@zju.edu.cn)

**Zhirui Zhang**

Tencent AI Lab

[zrustc11@gmail.com](mailto:zrustc11@gmail.com)

**Jun Xie**

Alibaba DAMO Academy

[qingjing.xj@alibaba-inc.com](mailto:qingjing.xj@alibaba-inc.com)

**Jordan Boyd-Graber**

CS, iSchool, UMIACS, LSC  
University of Maryland

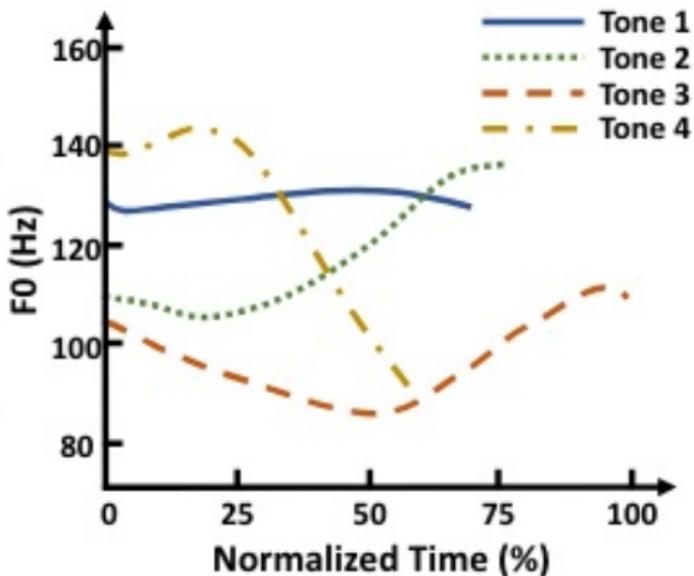
[jbgr@umiacs.umd.edu](mailto:jbgr@umiacs.umd.edu)

Tone 1: 叔 shū (uncle)

Tone 2: 熟 shú (cooked, familiar)

Tone 3: 鼠 shǔ (mouse, Muroidea)

Tone 4: 树 shù (tree)



Tones in Chinese (for “shu”, not “ma” like I said)

Original Lyrics  
(Inconsistent Tone)



sì → zài yǎn qián  
似 在 眼 前  
appear where eye front

As if before my eyes

Inter-syllable pitch alignment score: 0.5

Misheard Lyrics  
(Consistent Tone)



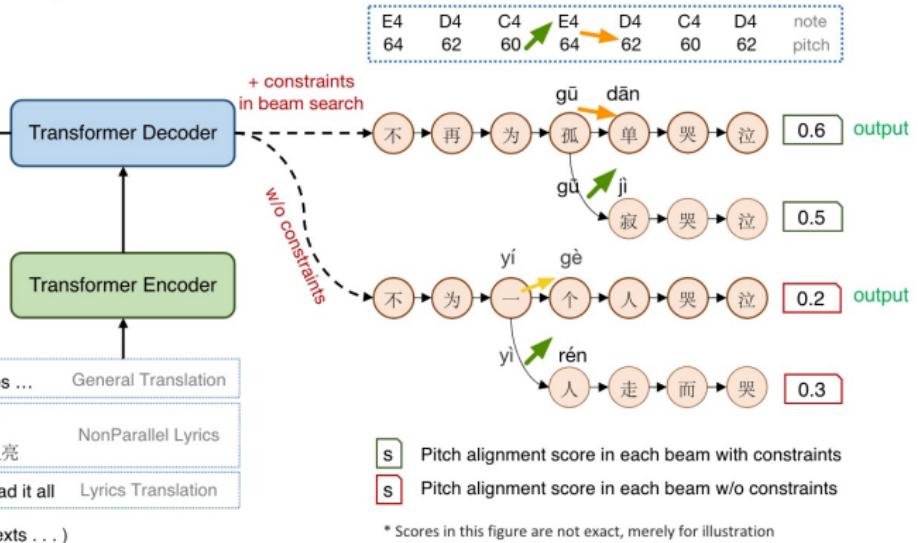
sǐ → zài yǎn qián  
死 在 眼 前  
death where eye front

Die before my eyes

Inter-syllable pitch alignment score: 0.75

Misheard lyrics when the tones are wrong

$y_4$  即便是美国…  
 $y_3$  Rolling in the deep  
 $y_2$  望向亘古无声的月亮  
 $y_1$  我们本来可以拥有一切



## Decoding song translations with tones in decoder

