# Do DALL-E and Flamingo Understand Each Other?

**Hang Li**[1,2]     **Jindong Gu**[3]     **Rajat Koner**[1]     **Sahand Sharifzadeh**[1]     **Volker Tresp**[1,2]
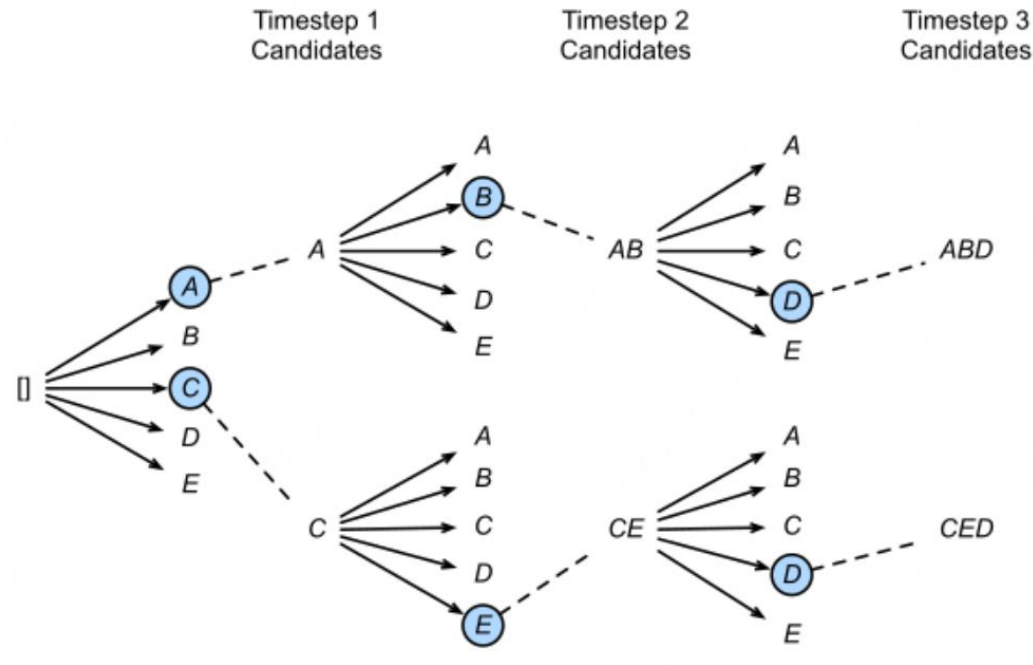
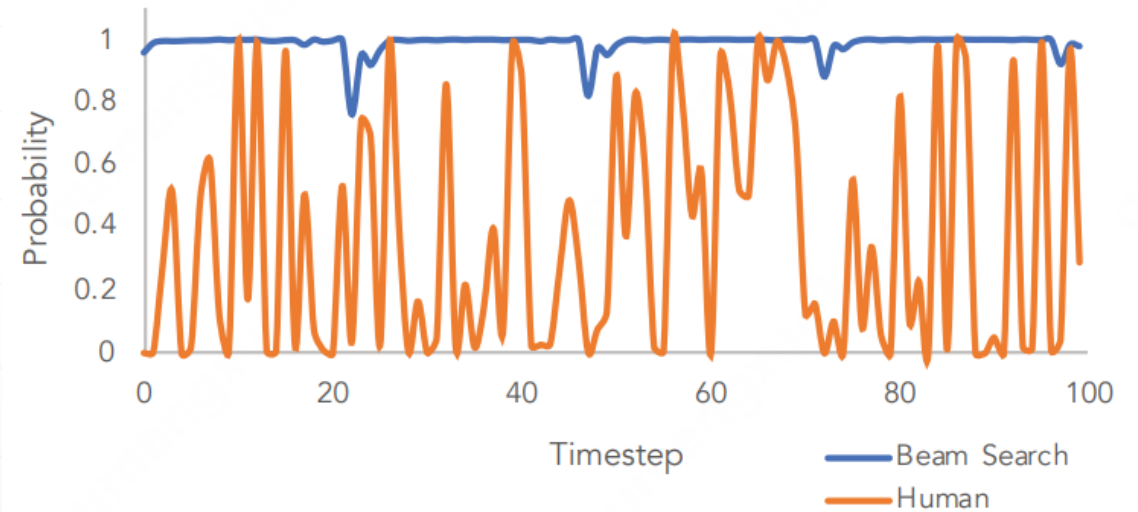[1]LMU Munich, Germany, [2]Siemens AG, Germany, [3]University of Oxford, UK
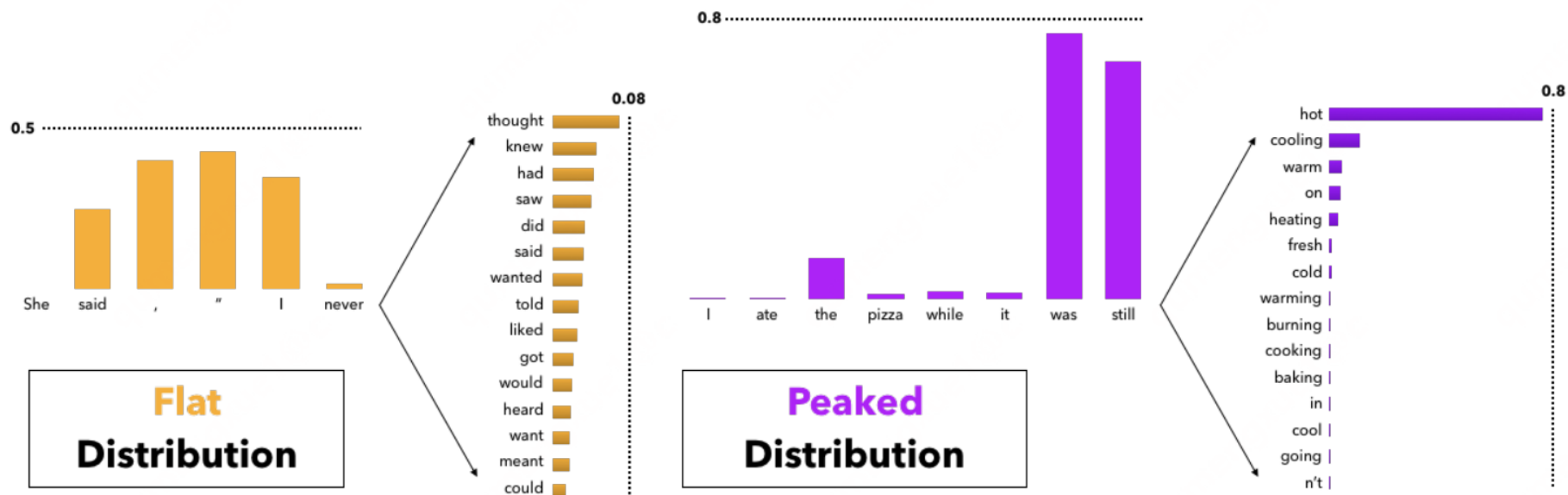
Mengxue

# Preview: Beam Search



Timestep 1 Candidates — Timestep 2 Candidates — Timestep 3 Candidates

Beam Search Text is Less Surprising

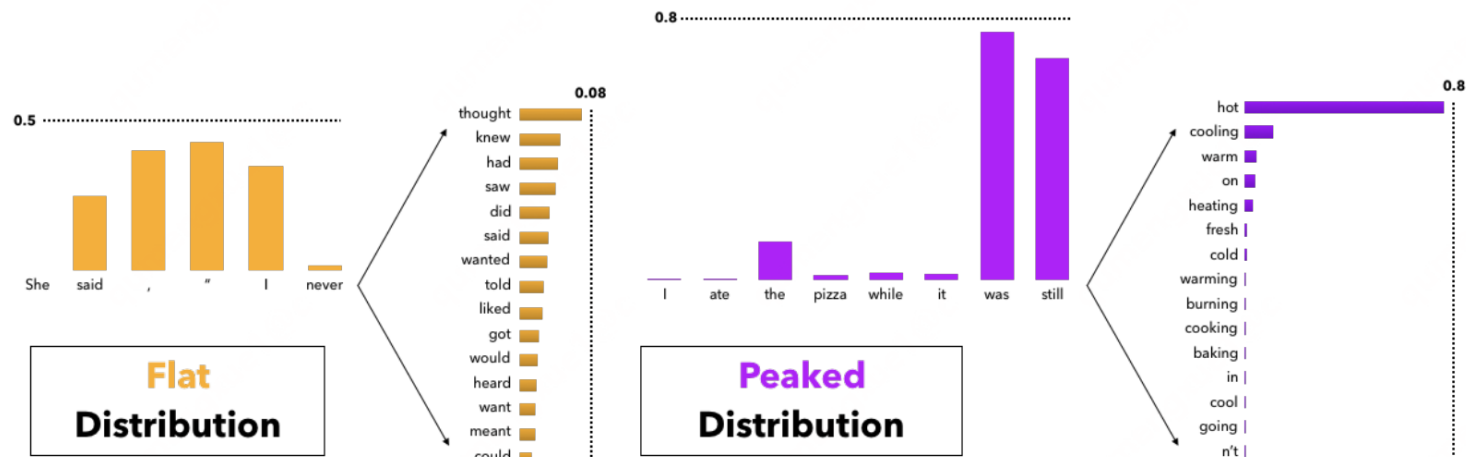Holtzman, Ari, et al. "The Curious Case of Neural Text Degeneration." *ICLR 2020*.

# Preview: Top-k Sampling



▪ 这个方法就是在采样前将输出的概率分布截断，取出概率最大的k个词构成一个集合，然后将这个子集词的概率再归一化，最后从新的概率分布中采样词汇。

▪ 但因为概率分布变化比较大，有时候可能很均匀(flat)，有的时候比较集中(peaked)。当分布均匀时，一个较小的k容易丢掉很多优质候选词。但如果k定的太大，这个方法会退化回全局随机采样。

Holtzman, Ari, et al. "The Curious Case of Neural Text Degeneration." *ICLR 2020*.
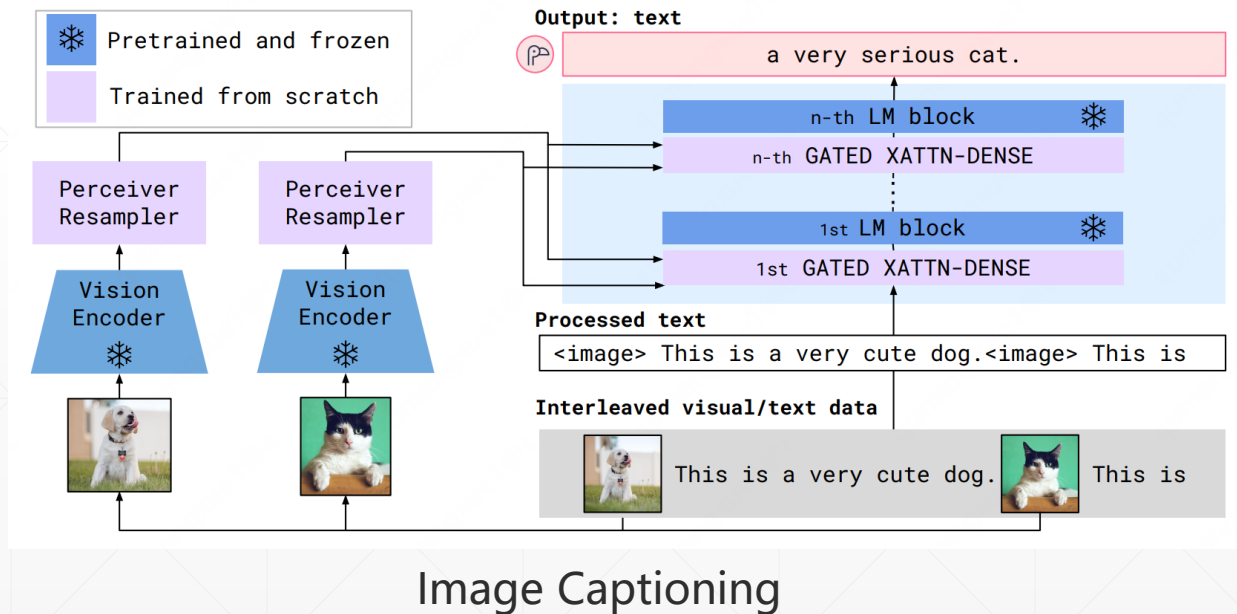
# Preview: Nucleus sampling(Top-p Sampling)



- 在每个时间步，解码词的概率分布满足80/20原则或者说长尾分布，头部的几个词的出现概率已经占据了绝大部分概率空间，把这部分核心词叫做nucleus。

- 基于这样的观察，提出nucleus sampling：给定一个概率阈值p，从解码词候选集中选择一个最小集Vp，使得它们出现的概率和大于等于p。然后再对Vp做一次re-scaling，本时间步仅从Vp集合中解码。

- 这样的好处在于在不同时间步，随着解码词的概率分布不同，候选词集合的大小会动态变化，不像top-k sampling是一个固定的窗口大小。由于解码词还是从头部候选集中筛选，这样的动态调整可以使生成的句子在满足多样性的同时又保持通顺。

Holtzman, Ari, et al. "The Curious Case of Neural Text Degeneration." *ICLR 2020.*

# What's the best choice?

- Beam Search → Top-k Sampling → Top-p Sampling

- **Caption** → Given an image, what textual description most accurately describes the image?

- **Generation** → Given a text, what is the best image that can present the semantics of the text?

**Do DALL-E and Flamingo Understand Each Other?**



Image Captioning

Text-to-Image Generation

(a) a tapir made of accordion. a tapir with the texture of an accordion.

(b) an illustration of a baby hedgehog in a christmas sweater walking a dog

**Flammingo**

**DALLE-E**

# Motivation

▪ In this work, we argue that the best text or caption for a given image is the text which would generate the image which is the most similar to that image.  ⟳ **Cycle Consistency**

▪ Likewise, the best image for a given text is the image that results in the caption which is best aligned with the original text.



Image Captioning                                                              Text-to-Image Generation

# Framework

# Experiments



Figure 3. Left: for each given image, the better the reconstructed image (shown in x-axis), the better the caption (shown in y-axis). Right: for each given text, the better the reconstructed text (shown in x-axis), the better the image (shown in y-axis).

| | NoCaps | | | | | | |
|---|---|---|---|---|---|---|---|
| Method | B1 | B2 | B3 | B4 | CIDEr | CIDErD | SPICE |
| Nucleus | 73.0 | 52.4 | 36.1 | 24.1 | 85.0 | 74.6 | 11.6 |
| Ours | **74.3** | **53.8** | **37.3** | **25.2** | **91.5** | **80.3** | **12.3** |
| Gain (%) | +1.8 | +2.7 | +3.5 | +4.2 | +7.6 | +7.7 | +6.3 |

| | COCO | | | | | | |
|---|---|---|---|---|---|---|---|
| Method | B1 | B2 | B3 | B4 | CIDEr | CIDErD | SPICE |
| Nucleus | **66.9** | 47.1 | 32.4 | 21.9 | 98.2 | 90.1 | 19.6 |
| Ours | **66.9** | **47.2** | **32.5** | **22.0** | **100.4** | **92.0** | **20.1** |
| Gain (%) | +0.0 | +0.1 | +0.4 | +0.3 | +2.2 | +2.1 | +2.2 |

# Qualitative Results

# Failure Cases



**Annotation Imperfection**                    **Generation Bias**

# Thanks