

# Fast Inference from Transformers via Speculative Decoding (ICML 2023)

---

Akira Kawata

# Abstract

- Problem
  - We need to run LLM model  $K$  times to get  $K$  tokens.
- Proposed method
  - Speculative decode
  - This enables us to run multiple token decode in parallel.
  - This can be applied to any model without any modification.
- Result
  - They achieved 2-3 times speed up.
  - After applied this method, the result token distribution exactly matches with the original one.

## 2. Speculative Decoding

---

## 2.1. Overview

- 構成要素
  - ターゲットモデル ( $M_p$ ) : 高速化したい大型モデル (高品質だが遅い)
  - 近似モデル ( $M_q$ ) : 下書き用の小型モデル (低コストで高速)
- プロセス
  - 下書き (Drafting) : 小型モデル  $M_q$  が  $\gamma$  個のトークンを先読み生成
  - 並列検証 (Parallel Evaluation) : 大型モデル  $M_p$  が、それらの候補を一度の実行で並列に検証
  - 確定と修正 (Correction) : 採用されたトークンを確定。最初の不採用箇所を  $M_p$  が修正し、新たなトークンとして出力
- メリット
  - 効率性 : 1回の重い計算 ( $M_p$ ) で、最大  $\gamma+1$  個のトークンを生成可能
  - 品質保証 : 出力される分布は、大型モデル単体の場合と完全に同一

## 2.2. Standardized Sampling

- argmax、top-k、nucleus、温度（temperature）調整など、一般に異なる処理として扱われる手法を、単一の共通フレームワークで扱う。
- あらゆるサンプリング手法は、最終的に調整された確率分布からの標準的なサンプリングとして定義できる。
- 具体例（argmax）：最大値以外の確率をゼロにして正規化した分布からのサンプリング、と読み替えることが可能。
- 本論文の以降の議論では、ターゲットモデル（ $M_p$ ）と近似モデル（ $M_q$ ）から得られる分布  $p(x)$  および  $q(x)$  は、指定の手法ですでに調整済みであると仮定する。
- これにより、個別のサンプリングアルゴリズムに依存しない、汎用的な推論高速化の議論が可能になる

## 2.3 Speculative Sampling

---

**Algorithm 1** SpeculativeDecodingStep

---

**Inputs:**  $M_p, M_q, prefix$ .

▷ **Sample  $\gamma$  guesses  $x_1, \dots, x_\gamma$  from  $M_q$  autoregressively.**

**for  $i = 1$  to  $\gamma$  do**

$q_i(x) \leftarrow M_q(prefix + [x_1, \dots, x_{i-1}])$

$x_i \sim q_i(x)$

**end for**

▷ **Run  $M_p$  in parallel.**

$p_1(x), \dots, p_{\gamma+1}(x) \leftarrow$

$M_p(prefix), \dots, M_p(prefix + [x_1, \dots, x_\gamma])$

▷ **Determine the number of accepted guesses  $n$ .**

$r_1 \sim U(0, 1), \dots, r_\gamma \sim U(0, 1)$

$n \leftarrow \min(\{i - 1 \mid 1 \leq i \leq \gamma, r_i > \frac{p_i(x)}{q_i(x)}\} \cup \{\gamma\})$

▷ **Adjust the distribution from  $M_p$  if needed.**

$p'(x) \leftarrow p_{n+1}(x)$

**if  $n < \gamma$  then**

$p'(x) \leftarrow \text{norm}(\max(0, p_{n+1}(x) - q_{n+1}(x)))$

**end if**

▷ **Return one token from  $M_p$ , and  $n$  tokens from  $M_q$ .**

$t \sim p'(x)$

**return**  $prefix + [x_1, \dots, x_n, t]$

---

- We can predict tokens from **M<sub>p</sub>** using tokens from **M<sub>q</sub>** exactly!!!

Example implementation: [Speculative-Sampling/speculative\\_sampling.py](#)

- [▷ Sample  \$\gamma\$  guesses  \$x\_1, \dots, x\_\gamma\$  from  \$M\_q\$  autoregressively.](#)
- [▷ Run  \$M\_p\$  in parallel.](#)
  - Input:  $[n\_batch, \gamma]$  / int32
  - output:  $[n\_batch, \gamma, n\_vocab]$  / float32
- [▷ Determine the number of accepted guesses  \$n\$ .](#)
  - Control flow changes depending on the calculation.
- [▷ Adjust the distribution from  \$M\_p\$  if needed.](#)
- [Sample from the adjusted distribution](#)
- [▷ Return one token from  \$M\_p\$ , and  \$n\$  tokens from  \$M\_q\$ .](#)

あとでじっくり見たほうがいいかも

# 3. Analysis

---

## 3.1. Number of Generated Tokens

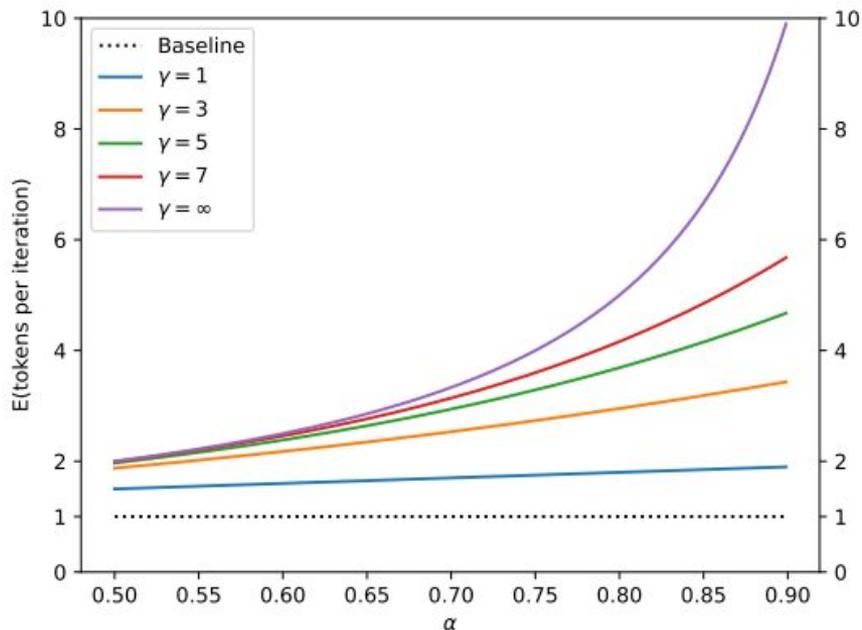


Figure 2. The expected number of tokens generated by Algorithm 1 as a function of  $\alpha$  for various values of  $\gamma$ .

- Vertical axis: The number of tokens generated per single speculative decoding step
- $\gamma$ : The number of tokens generated by the draft model.
- $\alpha$ : The degree of alignment between the probability distributions of  $M_q$  and  $M_p$
- So, red line should be close to 7 at  $\alpha = 1.0$ .

## 3.2 Calculating $\alpha$

飛ばす

### 3.3. Walftime Improvement

- $\alpha$  : 近似モデルが生成したトークンをターゲットモデルが採用する平均的な確率 (1.0 に近いほうがよい)
- $c$  : ターゲットモデルの 1 回の実行時間に対する、近似モデルの実行時間の比率 (小さいほうがよい)
- $T$  : Speculative decode なしの 1 token の生成時間
- $\gamma$  : Speculative に生成するトークン数 ( $\alpha$  と  $c$  に依存して最適値が決定)
- Speculative decode ありでの 1 token の生成時間は以下の式になる

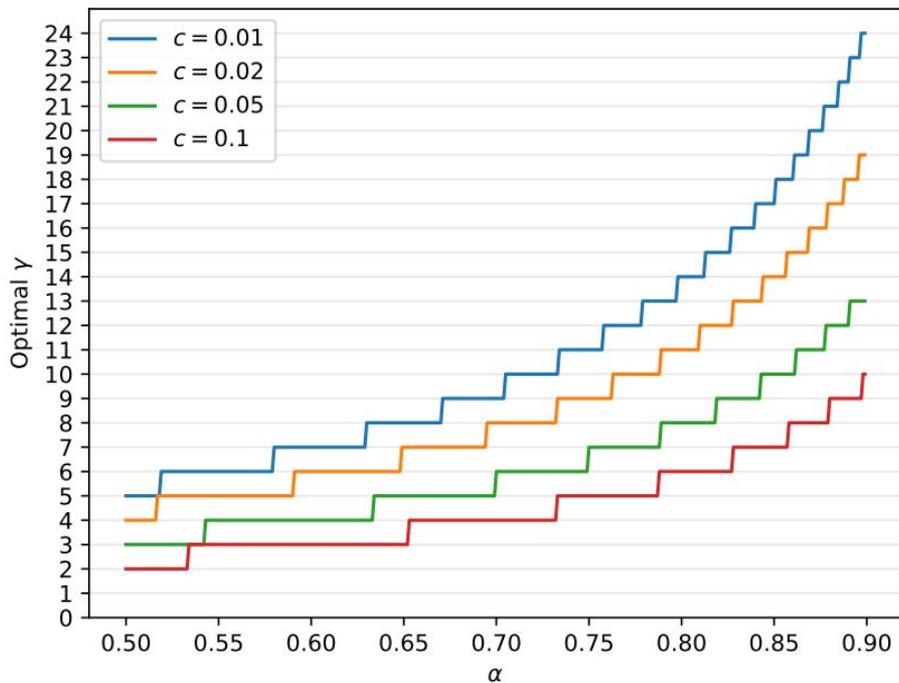
$$\frac{(c\gamma + 1)(1 - \alpha)}{1 - \alpha\gamma + 1} T$$

## 3.4. Number of Arithmetic Operations

- 演算量は以下の式に従って増える
  - $M_p$  の並列実行部分が無駄になるケースがあるので
- $\hat{c}$  は近似モデルとターゲットモデルの演算量の比
- 例えば  $\alpha = 0$  (ドラフトモデルがまったく当たらない)、 $\hat{c} = 0$  (ドラフトモデルの計算量が0) を入れると、計算量は  $\gamma + 1$  倍になる

$$\frac{(1 - \alpha)(\gamma \hat{c} + \gamma + 1)}{1 - \alpha\gamma + 1}$$

### 3.5. Choosing $\gamma$ (= どれだけ予測すればよいのか)



- ドラフトモデルの性能がよいほど (=  $\alpha$ が高い)、ドラフトモデルの計算コストが安いほど ( $c$ が小さい)、たくさん予測した法がよい

Figure 3. The optimal  $\gamma$  as a function of  $\alpha$  for various values of  $c$ .

## 3.6. Approximation Models

- Speculative decode はどんなドラフトモデルを使ってもよい!
- ドラフトモデルの候補
  - 二桁小さい Transformer モデル
  - n-gram モデル (テーブル引き)
  - コンテキストからトークンをコピーするモデル
  - 非自己回帰モデル
  - ランダム

# 4. Experiments

---

## 4.1. Empirical Walitime Improvement

Table 2. Empirical results for speeding up inference from a T5-XXL 11B model.

TASK	$M_q$	TEMP	$\gamma$	$\alpha$	SPEED
ENDE	T5-SMALL ★	0	7	0.75	<b>3.4X</b>
ENDE	T5-BASE	0	7	0.8	2.8X
ENDE	T5-LARGE	0	7	0.82	1.7X
ENDE	T5-SMALL ★	1	7	0.62	<b>2.6X</b>
ENDE	T5-BASE	1	5	0.68	2.4X
ENDE	T5-LARGE	1	3	0.71	1.4X
CNNNDM	T5-SMALL ★	0	5	0.65	<b>3.1X</b>
CNNNDM	T5-BASE	0	5	0.73	3.0X
CNNNDM	T5-LARGE	0	3	0.74	2.2X
CNNNDM	T5-SMALL ★	1	5	0.53	<b>2.3X</b>
CNNNDM	T5-BASE	1	3	0.55	2.2X
CNNNDM	T5-LARGE	1	3	0.56	1.7X

- ドラフトモデルとしては1/100ぐらいのモデルがよい
  - T5-XXL (11B)
  - T5-Large (800M)
  - T5-Base (250M)
  - T5-Small (77M)

## 4.2. Empirical $\alpha$ Values

Table 3. Empirical  $\alpha$  values for various target models  $M_p$ , approximation models  $M_q$ , and sampling settings. T=0 and T=1 denote argmax and standard sampling respectively<sup>6</sup>.

$M_p$	$M_q$	SMPL	$\alpha$
GPT-LIKE (97M)	UNIGRAM	T=0	0.03
GPT-LIKE (97M)	BIGRAM	T=0	0.05
GPT-LIKE (97M)	GPT-LIKE (6M)	T=0	0.88
GPT-LIKE (97M)	UNIGRAM	T=1	0.03
GPT-LIKE (97M)	BIGRAM	T=1	0.05
GPT-LIKE (97M)	GPT-LIKE (6M)	T=1	0.89
<hr/>			
T5-XXL (ENDE)	UNIGRAM	T=0	0.08
T5-XXL (ENDE)	BIGRAM	T=0	0.20
T5-XXL (ENDE)	T5-SMALL	T=0	0.75
T5-XXL (ENDE)	T5-BASE	T=0	0.80
T5-XXL (ENDE)	T5-LARGE	T=0	0.82
T5-XXL (ENDE)	UNIGRAM	T=1	0.07
T5-XXL (ENDE)	BIGRAM	T=1	0.19
T5-XXL (ENDE)	T5-SMALL	T=1	0.62
T5-XXL (ENDE)	T5-BASE	T=1	0.68
T5-XXL (ENDE)	T5-LARGE	T=1	0.71
<hr/>			
T5-XXL (CNNDM)	UNIGRAM	T=0	0.13
T5-XXL (CNNDM)	BIGRAM	T=0	0.23
T5-XXL (CNNDM)	T5-SMALL	T=0	0.65
T5-XXL (CNNDM)	T5-BASE	T=0	0.73
T5-XXL (CNNDM)	T5-LARGE	T=0	0.74
T5-XXL (CNNDM)	UNIGRAM	T=1	0.08
T5-XXL (CNNDM)	BIGRAM	T=1	0.16
T5-XXL (CNNDM)	T5-SMALL	T=1	0.53
T5-XXL (CNNDM)	T5-BASE	T=1	0.55
T5-XXL (CNNDM)	T5-LARGE	T=1	0.56
<hr/>			
LAMDA (137B)	LAMDA (100M)	T=0	0.61
LAMDA (137B)	LAMDA (2B)	T=0	0.71
LAMDA (137B)	LAMDA (8B)	T=0	0.75
LAMDA (137B)	LAMDA (100M)	T=1	0.57
LAMDA (137B)	LAMDA (2B)	T=1	0.71
LAMDA (137B)	LAMDA (8B)	T=1	0.74

- 様々なタスク、サンプリング手法、モデル構成において、理論上の期待採択率 ( $\alpha$ ) が実際にどの程度の値になるか調査した
- 超大型モデルにおいても、1/100 ぐらいのモデルで高い採択率を維持できる

## 5. Related work

---

## 5. Related work

SKIP!

# 6. Discussion

---

## 6. Discussion

SKIP!