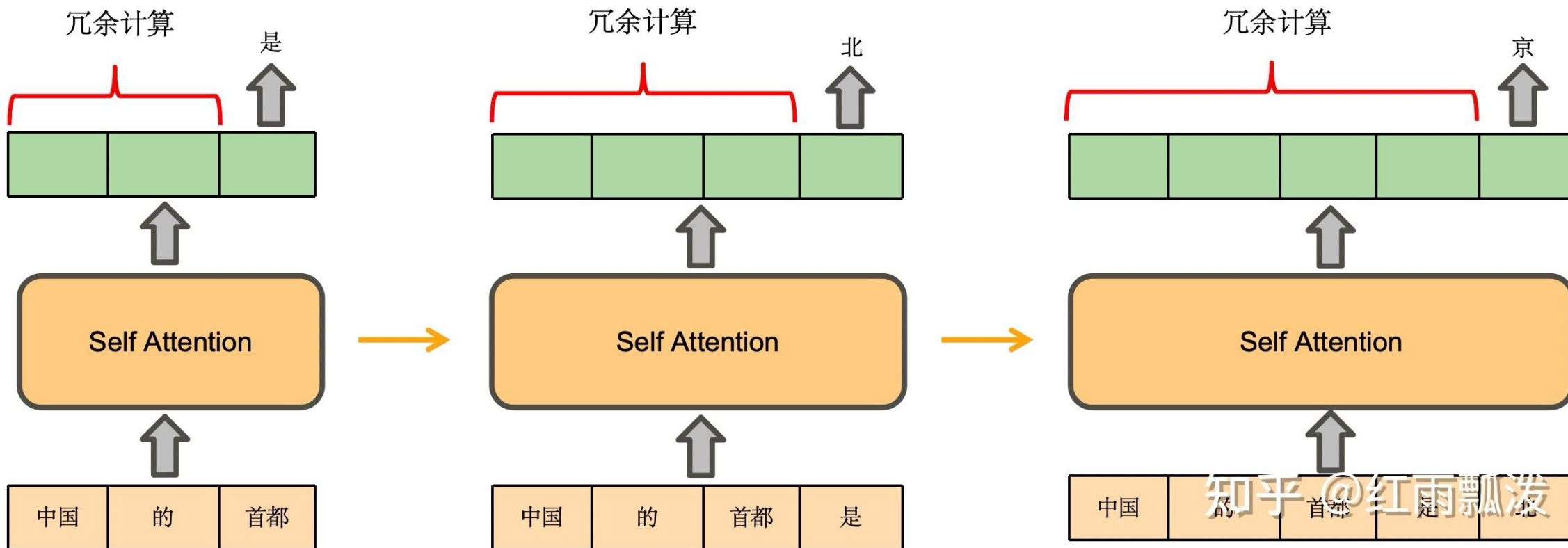# Speculative Decoding

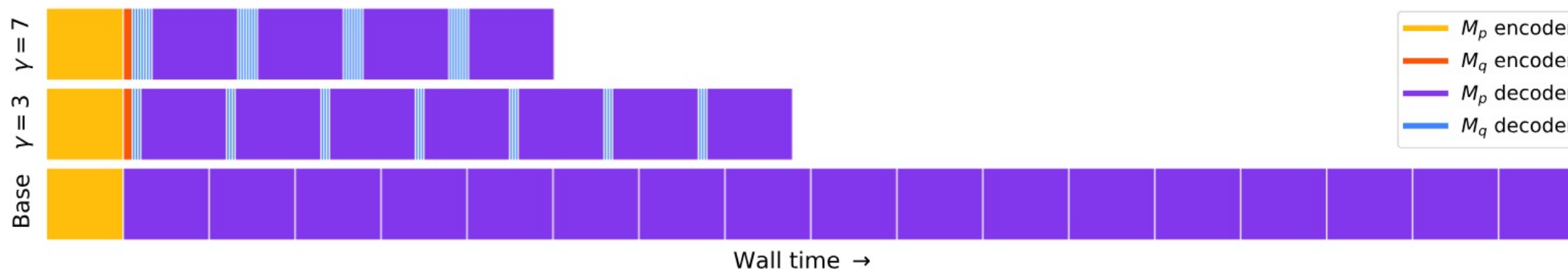# K-V Cache

# Speculative Decoding



title={Fast Inference from Transformers via Speculative Decoding},
arXiv={2211.17192},
citations={152}

# Speculative Decoding

$$q(x) \leq p(x) \quad \text{OK}$$

Random Sample

$$q(x) > p(x) \quad 1 - \frac{p(x)}{q(x)} \quad \text{If reject: resample} \quad p'(x) = norm(max(0, p(x) - q(x)))$$

Proof:

Note that as $p'(x) = norm(max(0, p(x) - q(x))) = \frac{p(x) - min(q(x), p(x))}{\sum_{x'}(p(x') - min(q(x'), p(x')))} = \frac{p(x) - min(q(x), p(x))}{1 - \beta}$,

$$P(x = x') = P(guess\ accepted, x = x') + P(guess\ rejected, x = x')$$

$$P(guess\ accepted, x = x') = q(x') min(1, \frac{p(x')}{q(x')}) = min(q(x'), p(x'))$$

$$P(guess\ rejected, x = x') = (1 - \beta)p'(x') = p(x') - min(q(x'), p(x'))$$

$$P(x = x') = min(p(x'), q(x')) + p(x') - min(p(x'), q(x')) = p(x').$$

# Speculative Decoding

$$E(\# \ generated \ tokens) = \frac{1 - \alpha^{\gamma+1}}{1 - \alpha}$$

$$\alpha = 1 - E(D_{LK}(p, q)) = E(\min(p, q))$$

| $\alpha$ | $\gamma$ | OPERATIONS | SPEED |
|---|---|---|---|
| 0.6 | 2 | 1.53X | 1.96X |
| 0.7 | 3 | 1.58X | 2.53X |
| 0.8 | 2 | 1.23X | 2.44X |
| 0.8 | 5 | 1.63X | 3.69X |
| 0.9 | 2 | 1.11X | 2.71X |
| 0.9 | 10 | 1.60X | 6.86X |

| TASK | $M_q$ | TEMP | $\gamma$ | $\alpha$ | SPEED |
|---|---|---|---|---|---|
| ENDE | T5-SMALL ★ | 0 | 7 | 0.75 | **3.4X** |
| 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 | **2.6X** |
| ENDE | T5-BASE | 1 | 5 | 0.68 | 2.4X |
| ENDE | T5-LARGE | 1 | 3 | 0.71 | 1.4X |
| CNNDM | T5-SMALL ★ | 0 | 5 | 0.65 | **3.1X** |
| CNNDM | T5-BASE | 0 | 5 | 0.73 | 3.0X |
| CNNDM | T5-LARGE | 0 | 3 | 0.74 | 2.2X |
| CNNDM | T5-SMALL ★ | 1 | 5 | 0.53 | **2.3X** |
| CNNDM | T5-BASE | 1 | 3 | 0.55 | 2.2X |
| CNNDM | T5-LARGE | 1 | 3 | 0.56 | 1.7X |

# Speculative Decoding

| $M_p$ | $M_q$ | 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 |
| 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 |
| 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 |
| 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 |

# Medusa

# Medusa

# Medusa



(a)



(b)

# Medusa

| Model Name | Vicuna-7B | Zephyr-7B | Vicuna-13B | Vicuna-33B |
|---|---|---|---|---|
| Acc. rate | 3.47 | 3.14 | 3.51 | 3.01 |
| Overhead | 1.22 | 1.18 | 1.23 | 1.27 |
| Quality | 6.18 (+0.01) | 7.25 (-0.07) | 6.43 (-0.14) | 7.18 (+0.05) |

# Thank You