
SpecTr: Fast Speculative Decoding via Optimal Transport

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Autoregressive sampling from large language models has shown to achieve state-of-
2 the-art results in several natural language tasks. However, autoregressive sampling
3 generates tokens one at a time making it slow, and even prohibitive in certain
4 tasks. One way to speed up decoding is *speculative decoding*: use a small model
5 to sample a *draft* (block or sequence of tokens), and then score all tokens in the
6 draft by the large language model in parallel. The tokens in the draft are either
7 accepted or rejected based on a statistical method to guarantee that the final output
8 is a valid sample from the large model. In this work, we provide a principled
9 understanding of speculative decoding through the lens of optimal transport (OT)
10 with *membership cost*. This framework can be viewed as an extension of the well-
11 known *maximal-coupling* problem. This new formulation enables us to generalize
12 the sampling method to allow for a set of k candidates at the token-level, which
13 leads to an improved optimal membership cost. We show that the optimal solution
14 can be computed via linear programming, whose best-known runtime is exponential
15 in k . We then propose an approximate solution whose acceptance probability is
16 $(1 - 1/e)$ -optimal multiplicatively. Moreover, it can be computed in time almost
17 linear with size of domain of a single token. Using this new OT algorithm, we
18 develop a new autoregressive sampling algorithm called *SpecTr*. We experimentally
19 demonstrate that the proposed approach achieves a speedup of 3X, a further 1.36X
20 speedup over speculative decoding on standard benchmarks.

21 1 Introduction

22 Autoregressive language models have shown to achieve state-of-the-art results in several natural
23 language tasks [2, 5, 20, 21]. During inference, given a context $x^t := x(1), x(2) \dots, x(t)$, an autore-
24 gressive model \mathcal{M}_b generates successive tokens $x(t+1), x(t+2), \dots$ via temperature sampling [1, 9],
25 where the next token $x(t+1)$ is drawn from the temperature-scaled distribution $\mathcal{M}_b(\cdot|x^t)$. If the
26 temperature is zero, i.e., greedy decoding, the next token is determined by the maximum likelihood
27 method i.e., $x(t+1) = \arg \max_{x \in \Omega} \mathcal{M}_b(x|x^t)$, where Ω is the vocabulary. The sampling approach
28 can be further combined with other sampling primitives such as nucleus sampling [13] and top- k
29 sampling [8, 17]. All these approaches are autoregressive decoding methods, where tokens are
30 generated serially one after another, which can be slow or even prohibitive in several applications
31 [18]. Hence, several techniques have been proposed to improve the speed of generation. Before we
32 proceed further, we first present some notations and a simplified computational model.

33 **Notations.** We use $x^{i:j}$ to denote the sequence $x(i), x(i+1), \dots, x(j)$ and when $i = 1$, we simply
34 use $x^j = x^{1:j}$. $x(i)$ denotes the i -th entry of x . Subscripts are used to distinguish between different
35 sequences. e.g., x_1^t and x_2^t denote two sequences of length t . We use $[n]$ to denote the set $\{1, \dots, n\}$.

36 **Computational model.**

- 37 • **Standard inference.** Given a context x^t , with $O(t^2)$ computation and $O(1)$ time, an
38 autoregressive model \mathcal{M}_b can compute $\mathcal{M}_b(y|x^t)$, the (temperature-scaled) probability of
39 all possible next tokens $y \in \Omega$.
- 40 • **Parallelization along the time axis.** Given a context x^t , with $O(t^2)$ computation and
41 $O(1)$ time, an autoregressive model \mathcal{M}_b can compute $\mathcal{M}_b(y|x^i)$, for all $y \in \Omega$ and $i \in$
42 $\{1, 2, \dots, t\}$.
- 43 • **Parallelization along time and batch axis¹.** Let K be the maximum batch size that
44 can be used during the inference of the autoregressive model. Given a several contexts,
45 $x_1^t, x_2^t, \dots, x_K^t$, with $O(Kt^2)$ computation and $O(1)$ time, an autoregressive model \mathcal{M}_b can
46 compute $\mathcal{M}_b(y|x_j^i)$, for all $y \in \Omega, i \in [t]$, and $j \in [K]$.

47 The above computation model shows that parallelizing along time and batch axes does not increase the
48 computation time. It is a simplified characterization of the typical hardware, such as TPUs and GPUs,
49 used in neural network inference. Previous approaches also assume similar computational model
50 to devise faster decoding algorithms [15, 4]. In practice, there will be some overhead depending
51 on hardware, implementation and resource utilization. In Section 8, we experimentally show that
52 the assumptions roughly hold for a large transformer model. We also note that there are efficient
53 transformer architectures, which reduces the computation cost from $O(t^2)$ to $O(t \log t)$ (see [19] for
54 a detailed survey). Such approaches are orthogonal to the focus of this paper, and they can be easily
55 combined with our approach.

56 Broadly speaking, multiple previous approaches proposed to guess a few possible future tokens using
57 an efficient model. They then compute several conditional probability distributions from the large
58 model based on the guesses. Computing the distributions takes $O(1)$ time due to parallelization. The
59 guessed tokens are then accepted or rejected based on a statistical method such that the accepted
60 tokens are effectively samples from the large model. When the guesses are good, multiple tokens will
61 be accepted. While this approach incurs the same computation cost as vanilla decoding (assuming
62 computing the guess is cheap), it can significantly improve decoding latency due to parallelization.

63 The goal of this work is to provide a principled understanding of the above approaches and discuss
64 optimality conditions and algorithmic improvements. We start by providing a more formal overview
65 of speculative decoding and related works.

66 **2 Previous works and speculative decoding**

67 Previous approaches make use of parallelization along the time axis to provide speed-ups. They first
68 predict multiple tokens and validate if these multiple tokens can be generated by the model with the
69 corresponding sampling or decoding scheme. For greedy decoding, multiple tokens can be predicted
70 by a separate model [18], aggressive decoding [10], or retrieval augmented text [23]. For sampling,
71 recently [15, 4] proposed an algorithm called speculative decoding, and we provide an overview of
72 this algorithm in the rest of the section.

73 Suppose we have access to a computationally-inexpensive draft model \mathcal{M}_s , which also predicts the
74 token given the context and the predictions of \mathcal{M}_s is similar to that of \mathcal{M}_b for most contexts. Suppose
75 we have decoded for t steps and have obtained prefix x^t . The next step of the speculative algorithm
76 can be broken down into three steps.

- 77 1. **Draft construction.** The draft model can be used to efficiently and “speculatively” sample
78 L tokens, extending the context to $x(1), x(2), \dots, x(t), \tilde{x}(t+1), \dots, \tilde{x}(t+L)$. We keep
79 the conditional probabilities on the next token $\mathcal{M}_s(y | x^t, \tilde{x}^{t+1:t+i})$ for each $i < L$ and
80 $\forall y \in \Omega$.
- 81 2. **Conditional probability computation.** After observing the samples, we compute the
82 conditional distributions $\mathcal{M}_b(y | x^t, \tilde{x}^{t+1:t+i})$ for each $i < L$ and $\forall y \in \Omega$ in parallel (along
83 time axis) in $O(1)$ time.
- 84 3. **Draft selection.** Select first L' of the L tokens and set $x(t+i) = \tilde{x}(t+i)$ for $i \leq L'$ given
85 the draft sequence and the conditional probabilities from both models.

¹This assumption also implies that naively batching multiple queries improves decoding throughput, but not latency.

86 After this step, we use $x_1^{t+L'}$ as prefix and sample the next sequence using speculative decoding
 87 iteratively. For completeness, we provide the full algorithm in Appendix A. The crux of the above
 88 three steps is draft selection, which given a draft sequence and the conditional probabilities from
 89 both models, selects a valid sequence such that the output has same distribution as that of the large
 90 model. In speculative decoding, this is achieved via recursively applying a token-level maximal
 91 coupling algorithm, which is provided in Algorithm 1. Note that for the draft selection, Algorithm 1 is
 92 applied where p is the conditional distribution of the draft model $\mathcal{M}_s(\cdot | x^t)$ and q is the conditional
 93 distribution of the large model $\mathcal{M}_b(\cdot | x^t)$ (which may be further conditioned on the context of the
 language model).

Algorithm 1 Token-level maximal coupling

Input: Distributions p, q , Draft sample $X \sim_{i.i.d.} p$.

- 1: Compute p^{res} where $\forall x \in \mathcal{X}, p^{\text{res}}(x) = \frac{q(x) - \min\{p(x), q(x)\}}{1 - \sum_{x'} \min\{p(x'), q(x')\}}$.
- 2: Set $Y = \perp$.
- 3: Sample $\eta \sim U(0, 1)$.
- 4: **if** $\eta \leq \min\left(1, \frac{q(X)}{p(X)}\right)$ **then**
- 5: $Y = X, \text{accept} = \text{True}$
- 6: **end if**
- 7: **Return** $Y \sim p^{\text{res}}, \text{accept} = \text{False}$.

94
 95 Algorithm 1 returns a random variable Y which either is the accepted input X ($\text{accept} = \text{True}$) or a
 96 sample from the residual distribution p^{res} ($\text{accept} = \text{False}$), which is defined in Step 1 of Algorithm 1.
 97 The algorithm is recursively applied as long as the draft tokens are accepted ($\text{accept} = \text{True}$) to
 98 select the first $L' \leq L$ tokens from the draft model. Previous works showed that if $X \sim p$, then $Y \sim q$
 99 [15, 4]. In the case of the draft selection this means that the output of the algorithm is distributed
 100 according to $\mathcal{M}_b(\cdot | x^t)$, which is exactly the desired outcome. Furthermore

$$\Pr(Y = X) = \sum_{x \in V} \min(p(x), q(x)) = 1 - d_{\text{TV}}(p, q),$$

101 where d_{TV} is the total variation distance between p and q . Since Y is distributed according to q , it is a
 102 valid sample from the large model. Secondly, the more similar p and q are, the higher the chance of
 103 $\Pr(Y = X)$, and fewer the number of serial calls to the larger model. In the ideal case, if $p = q$, then
 104 $\Pr(Y = X) = 1$, i.e., the draft token is always accepted, and when used for speculative decoding
 105 we have $L' = L$. In such a case, based on our computational model (Section 1), assuming the draft
 106 model is very fast compared to the large model, the speedup is L times.

107 3 Our contributions

108 From a theoretical viewpoint, the speculative decoding algorithm raises multiple questions.

- 109 • What is the relationship between speculative decoding and the broader literature of sampling in
 110 statistics?
- 111 • Is speculative decoding optimal in an information-theoretic sense?
- 112 • Speculative decoding uses parallelization along time to speed up decoding, would it be possible
 113 to use parallelization along batch (number of drafts) to further improve decoding speed?

114 We provide answers to all the above questions in this work. We first relate the problem of speculative
 115 decoding to the broader and well-studied discrete optimal transport theory (Section 4). With this
 116 connection, it becomes clear that the token-level draft selection is the optimal solution for optimal
 117 transport with indicator cost function and also related to the problem of maximal coupling [7]. Based
 118 on the connection to optimal transport, we show that one can further speed up the decoding by
 119 parallelizing along the batch axis by using multiple drafts from the draft model (Section 5).

120 More precisely, we formulate the draft selection problem as an discrete optimal transport problem
 121 with membership cost. Discrete optimal transport can be solved with a linear program, but the number
 122 of variables is exponential in batch size, which can be prohibitive. To address this, we propose an

123 approximate solution which achieves a $(1 - 1/e)$ -approximation of the optimal acceptance probability
 124 (Section 6).

125 With this theoretically motivated algorithm and guarantees, we circle back to speeding up decoding
 126 and propose a new algorithm called SpecTr and theoretically show that it can be used to derive valid
 127 sequences from the large model (Section 7). We then experimentally demonstrate the benefit of our
 128 approach on standard datasets (Section 8).

129 4 Token-level draft selection as an optimal transport problem

130 In this section, we formulate token-level draft as an optimal transport problem, where a cost function
 131 is associated with whether a draft token is accepted. To simplify notations, we assume the data comes
 132 from a discrete domain, but this can be easily generalized.

133 **Definition 1** (Coupling). For two probability distributions P over \mathcal{X} and Q over \mathcal{Y} , we say a joint
 134 distribution π supported over $\mathcal{X} \times \mathcal{Y}$ is a coupling between P and Q if

$$\begin{aligned} \forall y \in \mathcal{Y}, \quad \sum_{x \in \mathcal{X}} \pi(x, y) &= Q(y), \\ \forall x \in \mathcal{X}, \quad \sum_{y \in \mathcal{Y}} \pi(x, y) &= P(x). \end{aligned}$$

135 We use $\Pi(P, Q)$ to denote the set of all possible couplings between P and Q .

136 When it is clear from context, we will overload notation and refer to the probabilistic mapping
 137 $f_\pi : \mathcal{X} \rightarrow \mathcal{Y}$ introduced by the conditional probability $\pi(y | x) := \pi(x, y) / P(x)$ as a coupling, which
 138 is also referred to the transport plan from P to Q [22].

139 **Definition 2** (Optimal Transport (OT) [22]). For a cost function $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, the *transportation*
 140 *cost* of a coupling is defined as:

$$C(\pi) = \mathbb{E}_{X, Y \sim \pi} [c(X, Y)].$$

141 The *optimal transport plan* is the coupling $\pi \in \Pi(P, Q)$ that minimizes the transportation cost.

142 With these definitions in place, we can see that with $\mathcal{X} = \mathcal{Y} = \Omega$, which is the alphabet of the tokens,
 143 we recover the speculative decoding with the cost function of *indicator cost*, which captures the
 144 resampling cost, defined below:

$$\forall x \in \Omega, y \in \Omega, \quad c(x, y) = \mathbb{1}\{y \neq x\}.$$

145 The transportation cost of the coupling will be

$$C(\pi) = \mathbb{E}_{X, Y \sim \pi} [\mathbb{1}\{Y \neq X\}] = \mathbb{P}_{X, Y \sim \pi}(Y \neq X).$$

146 This optimal transport with this specific cost function is also called maximal coupling [7], and the
 147 optimal cost is known to be

$$\min_{\pi: \Pi} \mathbb{P}_{X, Y \sim \pi}(Y \neq X) = \sum_{x \in \Omega} \min(P(x), Q(x)). \quad (1)$$

148 Moreover, it can be shown that Algorithm 1 is equivalent to the maximal coupling between p and q ,
 149 and hence it achieves the optimal cost [7].

150 5 Optimal transport with multiple draft tokens

151 In this section, we generalize speculative decoding to allow for multiple drafts. More formally, let
 152 $\mathcal{X} = \Omega^k$ for some $k \in \mathbb{N}_+$, which is the set of k draft tokens from Ω and $\mathcal{Y} = \Omega$, which is the space
 153 of the final sampled token from the desired distribution. To characterize the resampling cost, we use
 154 the cost function of *membership cost*, defined below:

$$\forall x \in \Omega^k, y \in \Omega, \quad c(x, y) = \mathbb{1}\{y \notin S(x)\},$$

155 where $S(x) = \{o \in \Omega \mid o \text{ appears in } x\}$ denotes the set of distinct elements in x . When $k = 1$, this
 156 recovers the indicator cost mentioned above. The transportation cost of the coupling will be

$$C(\pi) = \mathbb{E}_{X,Y \sim \pi} [\mathbb{1}\{Y \notin S(X)\}] = \mathbb{P}_{X,Y \sim \pi}(Y \notin S(X)). \quad (2)$$

157 We will also refer to the above cost $C(\pi)$ as the *rejection probability* due to its probabilistic interpre-
 158 tation. And similarly, $\alpha(\pi) := 1 - C(\pi) = \mathbb{P}(Y \in S(X))$ will be the *acceptance probability*.

159 From now on we will use membership cost as the default cost function and refer to the optimal
 160 transport solution as *optimal transport with membership cost* (OTM). We use π^* to denote the
 161 coupling that minimizes this cost $\pi^* = \arg \min_{\pi \in \Pi(P,Q)} C(\pi)$; ² and the cost $C(\pi^*)$ is referred
 162 to as the *optimal transport cost* between P and Q . We use $\alpha(P, Q) = 1 - C(\pi^*)$ to denote the
 163 corresponding optimal acceptance probability.

164 **Draft selection with *i.i.d.* draft tokens.** In this paper, we will mainly focus on the case when the
 165 draft tokens are *i.i.d.* samples from a base distribution. Let p, q be supported over Ω and the goal is to
 166 obtain one valid token from q given k *i.i.d.* samples from p . This applies to the practical scenario
 167 where there exists a computationally efficient model, from which we can sample multiple independent
 168 draft tokens efficiently. We set $P = p^{\otimes k}$, a product distribution whose marginals are all p , and $Q = q$.
 169 And the OT problem we want to solve is the following:

$$\min C(\pi) \quad \text{s.t.} \quad \pi \in \Pi(p^{\otimes k}, q). \quad (3)$$

170 In this case, we overload notation and denote the *optimal acceptance probability* as
 171 $\alpha_k(p, q) := \alpha(p^{\otimes k}, q) = 1 - C(\pi^*)$. To better understand the quantity, below we state a few properties
 172 of α_k .

173 **Lemma 1.** (Appendix B.1) *The optimal acceptance probability satisfies the following properties.*

- 174 • **Monotonicity.** For any p, q and $k \geq 1$, $\alpha_k(p, q) \leq \alpha_{k+1}(p, q)$.
- 175 • **Consistency.** If $q(x)/p(x)$ is bounded $\forall x \in \Omega$, we have

$$\lim_{k \rightarrow \infty} \alpha_k(p, q) = 1.$$

176 Else,

$$\lim_{k \rightarrow \infty} \alpha_k(p, q) = \sum_{x \in \Omega} \mathbb{1}\{p(x) > 0\} q(x).$$

177 With the above result, it is clear that increasing k might increase the acceptance probability, particu-
 178 larly when the draft model satisfies $p(x) > 0$ for all $x \in \Omega$. We now focus on computing the optimal
 179 transport scheme and the optimal acceptance probability. Optimal transport in discrete domain has
 180 been studied extensively [14, 16, 11], and it is shown that the optimal transport problem is equivalent
 181 to the following linear programming problem:

$$\begin{aligned} \min & \sum_{x \in \Omega^k} \sum_{y \in \Omega} \pi(x, y) \mathbb{1}\{y \notin S(x)\} & (4) \\ \text{s.t.} & \quad \forall y \in \Omega, \sum_x \pi(x, y) = Q(y) \\ & \quad \forall x \in \Omega^k, \sum_y \pi(x, y) = P(x) \\ & \quad \forall x \in \Omega^k, y \in \Omega, \pi(x, y) \geq 0. \end{aligned}$$

182 Linear programming can be solved in time polynomial in the number of variables and constraints
 183 [6, 16]. Linear program in (4) has $|\Omega|^{k+1}$ variables and $|\Omega|^k + |\Omega|$ equality constraints.

184 **Lemma 2.** *Given p, q over Ω , there exists an algorithm that computes a solution to Eq. (3) in time*
 185 *$O(|\Omega|^{O(k)})$.*

²The existence of optimal coupling in discrete domain is well-known, e.g., see [22]. When the optimal coupling is not unique, we use π^* to denote one of the optimal couplings.

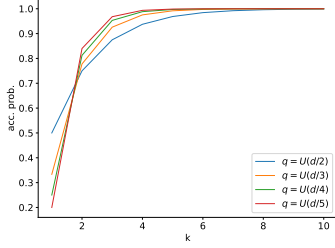


Figure 1: Optimal acc. prob. as a function of k when $p = U(d)$ for $d = 120$.

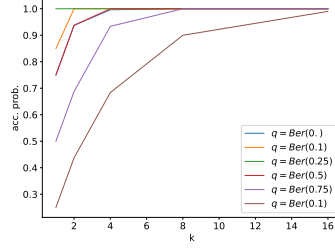


Figure 2: Optimal acc. prob. as a function of k when $p = \text{Ber}(0.25)$.

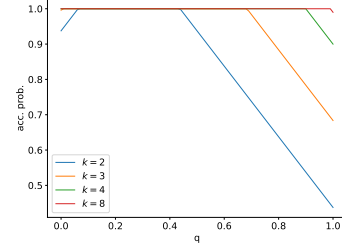


Figure 3: Optimal acc. prob. as a function of q when $p = \text{Ber}(0.25)$.

186 We refer to the optimal coupling obtained above as OTM- k and denote it as $\pi^{\text{OTM}-k}$. For the case of
 187 $k = 1$, we have a closed form expression for the optimal acceptance cost (see Eq. (1)), whereas for
 188 larger values of k , we don't have a general closed form expression.

189 We illustrate for few simple cases and plot them in Figures 1, 2, 3 and provide analysis for these
 190 simple distributions in Appendix B.2. Let $U(d)$ denote a uniform distribution over $[d]$. In Figure 1,
 191 we plot the optimal acceptance probability for different uniform functions q as a function of k .
 192 Observe that all acceptance probabilities are monotonically increasing and tend to one as $k \rightarrow \infty$,
 193 however the rate of convergence is vastly different. Furthermore if $\alpha_1(p, q) > \alpha_1(p, q')$, that does
 194 not necessarily mean $\alpha_k(p, q) > \alpha_k(p, q')$. In Figure 2, we plot the optimal acceptance probability
 195 for different Bernoulli distributions q as a function of k when $p = \text{Ber}(0.25)$. Note that when $p = q$,
 196 the acceptance probability is always one (green line), but as we increase / decrease q the acceptance
 197 probability decreases. Finally, in Figure 3, we plot the acceptance probability for different values of
 198 k as a function of q , when $p = \text{Ber}(0.25)$. In this scenario, note that if k is sufficiently large, say 8,
 199 then for most values of q , the acceptance probability is one, however if k is small, then the acceptance
 200 probability depends on how close p and q are. Even though, we don't have a closed form solution
 201 for general k , we provide an information-theoretic upper bound in the next theorem. For the case of
 202 $k = 1$, this upper bound matches the optimal acceptance probability of previous results. We also note
 203 that this bound is tight for all of the above examples.

204 **Theorem 1** (Appendix B.3). *For any two distributions p, q and $\forall k \geq 1$, we have*

$$\alpha_k(p, q) \leq \min_{\Omega_0 \subset \Omega} \left\{ \sum_{y \in \Omega_0} \min \{q(y), 1 - (1 - p(y))^k\} + \sum_{x^k \in \Omega^k} \min \left\{ \prod_{i=1}^k p(x_i), \sum_{y \in x^k \cap \Omega_0^c} q(y) \right\} \right\}.$$

205 While this solution gives the optimal transportation cost, if we aim to use generic linear program
 206 solver to solve (4), to the best of our knowledge, the best-known runtime will be exponential in k ,
 207 which can be prohibitive when either the vocabulary size Ω or the number of draft tokens k is large.
 208 In the next section, we will present an approximate solution to the OTM problem and show that for
 209 any pair of distributions, it gives a $(1 - 1/e)$ approximation to the *optimal acceptance probability*
 210 α_k .

211 6 Approximate OTM via k -sequential selection

212 In this section, we present sequential selection algorithm (K-SEQ), an approximate solution to the
 213 optimal transport problem in Eq. (3), which can be efficiently computed in time almost linear in $|\Omega|$
 214 and logarithmic in k . The algorithm is presented in Algorithm 2.

215 At a high-level, the algorithm goes over all k samples X_1, \dots, X_k generated from p sequentially, and
 216 decides on whether to accept each X_i based on the ratio $q(X_i)/p(X_i)$. The algorithm output the first
 217 accepted sample or result from a residual distribution p^{res} if none of the samples is accepted. To control
 218 the probability of accepting an $x \in \Omega$ with probability larger than $q(x)$. We choose an appropriate $\gamma \in$
 219 $[1, k]$ and accept X_i with probability $\min(1, q(X_i)/(\gamma \cdot p(X_i)))$ instead of $\min(1, q(X_i)/(p(X_i)))$
 220 as in the single-draft case. Further, notice that Algorithm 2 recovers Algorithm 1 when $\gamma = k = 1$.

Algorithm 2 k -sequential selection algorithm (K-SEQ).

Input: Distributions p, q , samples $X_1, \dots, X_k \sim i.i.d. p$. $\gamma \in [1, k]$: division factor.

1: Let $\beta_{p,q}(\gamma) = \sum_{x \in \Omega} \min(p(x), q(x)/\gamma)$ and $p_{\text{acc}} = 1 - (1 - \beta_{p,q}(\gamma))^k$. Compute p^{res} where

$$\forall x \in \Omega, p^{\text{res}}(x) = \frac{q(x) - \min\left\{p(x), \frac{q(x)}{\gamma}\right\} \frac{p_{\text{acc}}}{\beta_{p,q}(\gamma)}}{1 - p_{\text{acc}}}. \quad (5)$$

2: **for** $i = 1, 2, \dots, k$ **do**
 3: Sample $\eta_i \sim U(0, 1)$.
 4: **if** $\eta_i \leq \min\left(1, \frac{q(X_i)}{\gamma p(X_i)}\right)$ **then**
 5: $Y = X_i$.
 6: **Return** $Y = X_i$.
 7: **end if**
 8: **end for**
 9: **Return** $Y \sim p^{\text{res}}$.

221 In Theorem 2, we show that family of joint distributions induced by Algorithm 2 is indeed valid
 222 transportation plans from $p^{\otimes k}$ to q . Moreover, to find the best transportation plan within the family,
 223 we only need to search over a single parameter γ , which reduces the computation cost significantly.
 224 We also show that searching over this sub-family of couplings won't decrease the optimal acceptance
 225 probability by a multiplicative constant. The performance of Algorithm 2 is stated in Theorem 2.

226 **Theorem 2.** Let $\beta_{p,q}(\gamma) = \sum_{x \in \Omega} \min(p(x), \frac{q(x)}{\gamma})$ and γ^* be the solution to the identity below.

$$1 - (1 - \beta_{p,q}(\gamma))^k = \gamma \beta_{p,q}(\gamma). \quad (6)$$

227 When $\gamma \geq \gamma^*$, the coupling $\pi_\gamma^{\text{K-SEQ}}$ introduced by Algorithm 2 is a valid transport plan from $p^{\otimes k}$ to q
 228 and

$$\alpha(\pi_\gamma^{\text{K-SEQ}}) \geq p_{\text{acc}} = 1 - (1 - \beta_{p,q}(\gamma))^k.$$

229 And when $\gamma = \gamma^*$, we have

$$\alpha(\pi_{\gamma^*}^{\text{K-SEQ}}) \geq (1 - e^{-1}) \alpha_k(p, q).$$

230 Moreover, γ^* can be computed in time $O(|\Omega| \log k)$.

231 Due to space constraints, we provide the proof in the appendix. To see why γ^* can be computed
 232 efficiently, we notice that the function $f(\gamma)$ defined below has a root in $[1, k]$. Moreover it is
 233 continuous and monotonically increasing when $\gamma \in [1, k]$:

$$f(\gamma) = 1 - (1 - \beta_{p,q}(\gamma))^k - \gamma \beta_{p,q}(\gamma).$$

234 Hence the solution to Eq. (6) can be efficiently computed using binary search over the set $[1, k]$.

235 In fact, although Theorem 2 only guarantees that Algorithm 2 can achieve an acceptance rate at
 236 least a $(1 - e^{-1})$ factor of the optimal acceptance rate, empirically we observe that the acceptance
 237 probabilities are much closer for certain distributions. For example, for all the examples listed in the
 238 previous section, the proposed algorithm is in fact optimal. We list few more comparisons in the
 239 appendix.

240 7 SpecTr: Application of OTM in autoregressive sampling

241 In this section, we describe how OTM can be used to speed up auto-regressive sampling, which we
 242 refer to as SpecTr sampling. Similar to speculative decoding, each step of SpecTr can be decomposed
 243 into three phases:

- 244 1. **Draft set construction.** Given context x^T , use the draft model sample a set of draft
 245 sequences with length L , denoted by $S = \{z^L \sim \mathcal{M}_s(\cdot | x^T)\}$. We keep the conditional
 246 probabilities $\mathcal{M}_s(y | x^t, z^i)$ for all $y \in \Omega, i \leq L$ and $z^L \in S$.
- 247 2. **Conditional probability computation.** Compute the conditional probabilities on the next
 248 token for the large model $\mathcal{M}_b(y | x^t, z^i)$ for all $y \in \Omega, i \leq L$ and $z^L \in S$ in parallel.

249
250

3. **Draft selection.** Select first L' of the L tokens and set $x(t+i) = z(i)$ for $i \leq L'$ and some $z \in S$ given the set of draft sequences and the conditional probabilities from both models.

Algorithm 3 Draft selection with multiple candidates (DraftSelection).

Input: Input sequence x^t ; candidate length: L ; a set of candidates $S = \{z_i^L \mid i = 1, \dots, |S|\}$ with length L .

- 1: Compute a transport plan (using linear programming in Lemma 2 for an exact solution or Algorithm 2 for an approximate solution) from $\mathcal{M}_s(\cdot \mid x^t)^{\otimes |S|}$ to $\mathcal{M}_b(\cdot \mid x^t)$, denoted by π_t .
- 2: Get the multi-set of next token-level drafts: $S_z = \{z_i(1)\}_{i \in [|S|]}$ and compute $Y = f_{\pi_t}(S_z)$.
- 3: **if** $L = 1$ **then**
- 4: **if** $Y \in S_z$ **then**
- 5: Sample $Y' \sim \mathcal{M}_b(\cdot \mid (x^t, Y))$
- 6: **Return** (x^t, Y, Y') .
- 7: **else**
- 8: **Return** (x^t, Y)
- 9: **end if**
- 10: **end if**
- 11: Let $S_{\text{next}} = \{z^{2:L} \mid z \in S \text{ and } z(1) = Y\}$ be the set that consists of sub-sequences of the candidates that agree with the selected next token.
- 12: **if** $S_{\text{next}} = \emptyset$ **then**
- 13: **Return** (x^t, Y)
- 14: **else**
- 15: **Return** DraftSelection($(x^t, Y), L - 1, S_{\text{next}}$)
- 16: **end if**

251 **Draft set with *i.i.d.* draft sequences.** Given
252 context x^t , a natural way to come up with a
253 set of K drafts is to independently sample K
254 draft sequences from the conditional distribution
255 $\mathcal{M}_s(\cdot \mid x^t)$, *i.e.*,

$$z_1^L, z_2^L, \dots, z_K^L \sim_{i.i.d.} \underbrace{\mathcal{M}_s(\cdot, \cdot, \dots \mid x^t)}_{L \text{ dots}} \quad (7)$$

256 The draft set construction method in (7) can
257 be generalized to a prefix-tree based algorithm.
258 However, this generalized version did not per-
259 form better in experiments. We include this
260 construction in the appendix for completeness

261 **Draft selection with multiple candidates.** We
262 present the selection algorithm given a set of
263 draft sequences in Algorithm 3. We assume
264 the condition probabilities on the next token is
265 available given any prefix in the candidate set
266 since they are computed parallelly in the second
267 phase, and won't list them as inputs explicitly
268 in Algorithm 3.

269 A sample run of the algorithm is presented in
270 Fig. 4. The algorithm proceeds in a recursive
271 fashion. Given prompt x^t and a candidate set S sampled from $\mathcal{M}_s(\cdot \mid x^t)$, the algorithm first
272 computes the optimal transport plan $f_\pi : \Omega^{|S|} \rightarrow \Omega$ from $\mathcal{M}_s(\cdot \mid x^t)^{\otimes |S|}$ to $\mathcal{M}_b(\cdot \mid x^t)$. Then f_π
273 is applied to the first token in each sequence in S to obtain a valid token Y from $\mathcal{M}_b(\cdot \mid x^t)$. If
274 Y is not the last token ($L \geq 2$), we filter out sequences in S whose first token is not Y and denote
275 the remaining sequences as S_{next} and feed it to the algorithm with context (x^t, Y) and draft length
276 $L - 1$. This goes on until we have $L = 1$ or $S_{\text{next}} = \emptyset$.

277 In this case when Y is the last token (*i.e.*, $L = 1$) and $Y \in S$, we have the choice to sample an
278 additional token $\mathcal{M}_b(\cdot \mid (x^t, Y))$ since this conditional probability is already computed in the second

$|S_z| = 6$ $|S_z| = 3$ $|S_z| = 2$ $|S_z| = 1$

be	liked	by	all
be	read	by	four
be	liked	for	its
not	be	liked	by
not	get	good	reviews
receive	one	good	review

This paper will

Figure 4: An example of draft selection in SpecTr with $L = 4$ and $K = 6$. Draft selection algorithm has input of all conditional probabilities from both large and small models. In the first step, we compute the transport plan with $|S_z| = K = 6$ and the sequential selection algorithm will select 'be', which appeared thrice in our samples. We then compute the transport plan with $|S_z| = 3$ and the sequential selection algorithm will select 'liked'. We then compute the transport plan with $|S_z| = 2$ and the sequential selection algorithm will select 'by'. Finally, we compute the transport plan with $|S_z| = 1$ and the sequential selection algorithm will not select any of the drafts.

Table 1: Average latency with parallelization along the time axis and batch axis. We report average latency with standard deviation over 1,000 runs using a 97M transformer relative to length = 4 and batch = 1 on GPU.

Relative latency	batch = 1	batch = 2	batch = 4	batch = 8
length = 4	1.00 ± 0.16	1.01 ± 0.15	1.06 ± 0.10	1.10 ± 0.16
length = 8	1.01 ± 0.18	1.09 ± 0.25	1.10 ± 0.09	1.42 ± 0.4

Table 2: Experimental results on the LM1B dataset. All results are over 1000 test prompts averaged over three different random seeds.

Algorithm	K	L	Number of decoded tokens per serial call
Baseline	-	-	1.0
Speculative	1	4	2.2
SpecTr	2	4	2.4
SpecTr	4	4	2.7
SpecTr	8	4	3.0
Speculative	1	8	2.3
SpecTr	2	8	2.6
SpecTr	4	8	3.0
SpecTr	8	8	3.3

279 phase. Due to the property of the transport plan, we know that Y is always a valid sample from
 280 $\mathcal{M}_b(\cdot | x^t)$. The overall performance of the algorithm is stated in Theorem 3. We needed to take care
 281 in the statement and the proof to deal with the fact that the length τ of the output sequence Y^τ is
 282 itself a random variable. We defer the proof to the appendix due to limited space.

283 **Theorem 3.** Assume all drafts in set S are generated from the small model with input x^t , or more
 284 precisely, $\forall z \in S$,

$$\forall i \in [1, L], \quad z(i) \sim \mathcal{M}_b(\cdot | x^t, z^{i-1}). \quad (8)$$

285 Let Y^τ be the output of Algorithm 3 where τ is the length of the output, and $Z^{\tau+1:L} = (Z(\tau +$
 286 $1), \dots, Z(L)) \sim \mathcal{M}_b(\underbrace{\cdot, \cdot, \dots}_{(L-\tau) \text{ dots}} | x^t, Y^\tau)$, then it satisfies that $(Y^\tau, Z^{\tau+1:L}) \sim_{\text{prob}} \mathcal{M}_b(\underbrace{\cdot, \cdot, \dots}_{L \text{ dots}} |$
 287 $x^t)$. More precisely, For any length- L sequence $o^L = (o(1), \dots, o(L)) \in \Omega^L$, we have

$$\Pr((Y^\tau, Z^{\tau+1:L}) = o^L) = \prod_{i=1}^L \mathcal{M}_b(o(i) | x^t, o^{i-1}).$$

288 8 Experiments

289 We evaluate the performance of our algorithm and compare it to speculative decoding by following
 290 a recipe provided in [15]. We train decoder-only transformer models on the one-billion language
 291 benchmark (LM1B) [3] based on the example provided in the FLAX library [12]. For the draft
 292 model, we use a 6M parameter transformer model, and for the large model we use a 97M parameter
 293 transformer model.

294 We first provide a verification of the computational model introduced in Section 1 by reporting the
 295 latencies of using the large model to compute the probabilistic distributions with parallelization over
 296 time and batch axes. As shown in Table 1, the latency stays roughly constant in these setting.

297 The results of different decoding algorithms are shown in Table 2. The baseline method decodes one
 298 token from the large model per serial call, and speculative decoding improves this to ≈ 2.3 . The
 299 proposed method SpecTr improves upon speculative decoding and increases the number of decoded
 300 tokens per serial call as we increase the number of drafts K . We further note that for both Speculative
 301 decoding and SpecTr, the number of decoded tokens increases as we increase the block length from 4
 302 to 8. We also note that based on our current implementation, generating the drafts using the small
 303 models adds about 10%-15% latency under settings in Table 2. Due to space constraints, we provide
 304 additional experiments and details in Appendix F.

References

- 305
- 306 [1] David H Ackley, Geoffrey E Hinton, and Terrence J Sejnowski. A learning algorithm for
307 boltzmann machines. *Cognitive science*, 9(1):147–169, 1985.
- 308 [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
309 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
310 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 311 [3] Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and
312 Tony Robinson. One billion word benchmark for measuring progress in statistical language
313 modeling. *arXiv preprint arXiv:1312.3005*, 2013.
- 314 [4] Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and
315 John Jumper. Accelerating large language model decoding with speculative sampling. *arXiv
316 preprint arXiv:2302.01318*, 2023.
- 317 [5] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
318 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:
319 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 320 [6] George B Dantzig. Linear programming. *Operations research*, 50(1):42–47, 2002.
- 321 [7] Frank Den Hollander. Probability theory: The coupling method. *Lecture notes available online
322 (<http://websites.math.leidenuniv.nl/probability/lecturenotes/CouplingLectures.pdf>)*, 2012.
- 323 [8] Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. *arXiv
324 preprint arXiv:1805.04833*, 2018.
- 325 [9] Jessica Fidler and Yoav Goldberg. Controlling linguistic style aspects in neural language
326 generation. *arXiv preprint arXiv:1707.02633*, 2017.
- 327 [10] Tao Ge, Heming Xia, Xin Sun, Si-Qing Chen, and Furu Wei. Lossless acceleration for seq2seq
328 generation with aggressive decoding. *arXiv preprint arXiv:2205.10350*, 2022.
- 329 [11] Wenshuo Guo, Nhat Ho, and Michael Jordan. Fast algorithms for computational optimal trans-
330 port and wasserstein barycenter. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings
331 of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume
332 108 of *Proceedings of Machine Learning Research*, pages 2088–2097. PMLR, 26–28 Aug 2020.
- 333 [12] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas
334 Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2023.
- 335 [13] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural
336 text degeneration. *arXiv preprint arXiv:1904.09751*, 2019.
- 337 [14] Leonid V Kantorovich. On the translocation of masses. In *Dokl. Akad. Nauk. USSR (NS)*,
338 volume 37, pages 199–201, 1942.
- 339 [15] Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via
340 speculative decoding. *arXiv preprint arXiv:2211.17192*, 2022.
- 341 [16] Ofir Pele and Michael Werman. Fast and robust earth mover’s distances. In *2009 IEEE 12th
342 international conference on computer vision*, pages 460–467. IEEE, 2009.
- 343 [17] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
344 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 345 [18] Mitchell Stern, Noam Shazeer, and Jakob Uszkoreit. Blockwise parallel decoding for deep
346 autoregressive models. *Advances in Neural Information Processing Systems*, 31, 2018.
- 347 [19] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey.
348 *ACM Computing Surveys*, 55(6):1–28, 2022.

- 349 [20] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-
350 Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for
351 dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.
- 352 [21] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-
353 thé Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open
354 and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 355 [22] Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2009.
- 356 [23] Nan Yang, Tao Ge, Liang Wang, Binxing Jiao, Daxin Jiang, Linjun Yang, Rangan Majumder,
357 and Furu Wei. Inference with reference: Lossless acceleration of large language models. *arXiv*
358 *preprint arXiv:2304.04487*, 2023.

Algorithm 4 Speculative Sampling SPECSAMPLE.

Input: Input sequence x^t . Access to a small model \mathcal{M}_s and a large model \mathcal{M}_b , block length L , end of sequence symbol eos .

- 1: Autoregressively sample \mathcal{M}_s with context x^t to get $L - 1$ subsequent samples denoted by $\tilde{x}_{t+1}, \dots, \tilde{x}_{t+L-1}$.
- 2: Let $\tilde{x}_i = x_i$ for $i \leq n$.
- 3: In parallel compute $p_i = \mathcal{M}_s(\cdot | \tilde{x}^{t+i-1})$ and $q_i = \mathcal{M}_b(\cdot | \tilde{x}^{t+i-1})$ for $1 \leq i \leq L$.
- 4: **for** $i = 1, \dots, L - 1$ **do**
- 5: Compute $Y_i, accept = \text{Algorithm 1}(p_i, q_i, \tilde{x}_{t+i})$
- 6: $x_{t+i} = Y_i$.
- 7: **if** $x_{t+i} = eos$ **then**
- 8: **Return** x^{t+i}
- 9: **end if**
- 10: **if** $accept = \text{True}$ **then**
- 11: **Continue.**
- 12: **else**
- 13: **Return** SPECSAMPLE($x^{t+i}, \mathcal{M}_s, \mathcal{M}_b, L$).
- 14: **end if**
- 15: **end for**
- 16: Draw x_{t+L} from q_L .
- 17: **Return** SPECSAMPLE($x^{t+L}, \mathcal{M}_s, \mathcal{M}_b, L$).

B Missing proofs in Section 5**B.1 Proof of Lemma 1**

362 We first prove *monotonicity*. By definition,

$$\begin{aligned} \alpha_k(p, q) &= 1 - \min_{\pi \in \Pi(p^{\otimes k}, q)} \Pr_{X^k, Y \sim \pi}(Y \notin S(X^k)) \\ &= \max_{\pi \in \Pi(p^{\otimes k}, q)} \Pr_{X^k, Y \sim \pi}(Y \in S(X^k)) \end{aligned}$$

363 Moreover, for any $\pi \in \Pi(p^{\otimes k}, q)$, we can construct $\pi' \in \Pi(p^{\otimes k+1}, q)$ by setting

$$\forall x^{k+1} \in \Omega^{k+1}, y \in \Omega, \pi'(x^{k+1}, y) = \pi(x^k, x(k+1), y)p(x(k+1)),$$

364 *i.e.*, adding an independent sample from p to X^k .

365 Hence we have

$$\begin{aligned} \alpha_k(p, q) &= \max_{\pi \in \Pi(p^{\otimes k}, q)} \Pr_{X^k, Y \sim \pi}(Y \in S(X^k)) \\ &= \max_{\pi \in \Pi(p^{\otimes k}, q)} \Pr_{X^{k+1}, Y \sim \pi'}(Y \in S(X^k)) \\ &\leq \max_{\pi \in \Pi(p^{\otimes k}, q)} \Pr_{X^{k+1}, Y \sim \pi'}(Y \in S(X^{k+1})) \\ &\leq \max_{\pi' \in \Pi(p^{\otimes k+1}, q)} \Pr_{X^{k+1}, Y \sim \pi'}(Y \in S(X^{k+1})) \\ &= \alpha_{k+1}(p, q). \end{aligned}$$

366 Next we prove *consistency*. We start with the case when $\forall x \in \Omega, q(x)/p(x) < \infty$. To prove this, we
 367 will show that Algorithm 2 with $\gamma_{\max} = \max_{x \in \Omega} q(x)/p(x)$ satisfies

$$\lim_{k \rightarrow \infty} \alpha(\pi_{\gamma_{\max}}^{\text{K-SEQ}}) = 1.$$

368 Notice that by Lemma 3 and Theorem 2, $\pi_{\gamma_{\max}}^{\text{K-SEQ}}$ is a valid coupling, and

$$\alpha(\pi_{\gamma_{\max}}^{\text{K-SEQ}}) = 1 - (1 - \beta_{p,q}(\gamma_{\max}))^k,$$

369 where $\beta_{p,q}(\gamma) = \sum_{x \in \Omega} \min(p(x), \frac{q(x)}{\gamma}) \geq 1/\gamma_{\max} > 0$. Taking $k \rightarrow \infty$ concludes the proof.

370 For the case when $q(x)/p(x)$ is unbounded, there exists $x \in \Omega$ such that $q(x) > 0$ and $p(x) = 0$. Let

$$p_{\text{off}} = \sum_{x \in \Omega} \mathbb{1}\{p(x) = 0\} q(x).$$

371 Let x_0 be such that $p(x_0) > 0$. We define q' such that

$$q' = \begin{cases} 0, & \text{if } p(x) = 0, \\ q(x), & \text{if } p(x) > 0 \text{ and } x \neq x_0, \\ q(x) + p_{\text{off}} & \text{if } x = x_0. \end{cases}$$

372 Then we have $d_{\text{TV}}(q, q') = p_{\text{off}}$, and hence by subadditivity of transport cost,

$$\alpha_k(p, q) \geq \alpha_k(p, q') - p_{\text{off}}.$$

373 Moreover, we have $\forall x \in \Omega, q'(x)/p(x) < \infty$. Hence

$$\lim_{k \rightarrow \infty} \alpha_k(p, q) \geq \lim_{k \rightarrow \infty} \alpha_k(p, q') - p_{\text{off}} = 1 - p_{\text{off}} = \sum_{x \in \Omega} \mathbb{1}\{p(x) > 0\} q(x).$$

374 B.2 Optimal acceptance probability calculations

375 In this section, we provide a sketch of optimal acceptance probability calculations for results in
376 Figures 1, 2, and 3.

377 **Figure 1:** $p = U(d)$ and $q = U(d/r)$. The optimal acceptance probability is

$$\alpha_k(U(d), U(d/r)) = 1 - (1 - 1/r)^k.$$

378 We first prove $\alpha^k(U(d), U(d/r)) \geq 1 - (1 - 1/r)^k$ by a construction. Let $S(X^k)$ be the set of
379 unique symbols in X^k . Consider the following transport plan, where Y is drawn uniformly from
380 $S(X^k) \cap [d/r]$ and draws a new uniform sample from $[d/r]$ if $S(X^k) \cap [d/r] = \emptyset$. Observe that
381 since $U(d)$ is uniform over $[d]$, this is a valid transport plan and furthermore,

$$\alpha_k(U(d), U(d/r)) \geq \Pr(S(X^k) \cap [d/r] \neq \emptyset) = 1 - (1 - 1/r)^k.$$

382 The upper bound follows by setting $\Omega_0 = [d] \setminus [d/r]$ in Theorem 1.

$$\alpha_k(U(d), U(d/r)) \leq \Pr(S(X^k) \cap [d/r] \neq \emptyset) = 1 - (1 - 1/r)^k.$$

383 **Figure 2 and 3: Ber(p) and Ber(q).** The optimal acceptance probability is

$$\alpha_k(\text{Ber}(p), \text{Ber}(q)) = \min(q, 1 - (1 - p)^k) + \min(1 - q, 1 - p^k).$$

384 Setting $\Omega_0 = \{0, 1\}$ in Theorem 1 yields the upper bound. For the lower bound observe that since
385 $\Omega = \{0, 1\}$, $\mathbb{1}\{y \notin S(x^k)\} < 1$ if and only if x^k is 0^k or 1^k . Hence,

$$\alpha_k(\text{Ber}(p), \text{Ber}(q)) = \pi(X^k \notin \{0^k, 1^k\}) + \max_{\pi} \{\pi(Y = 0, X^k = 0^k) + \pi(Y = 1, X^k = 1^k)\}.$$

386 Consider the transport plan π given by $\pi(1^k, 1) = \min(p^k, q)$, $\pi(1^k, 0) = p^k - \min(p^k, q)$,
387 $\pi(0^k, 0) = \min((1 - p)^k, 1 - q)$, and $\pi(0^k, 1) = (1 - p)^k - \min((1 - p)^k, 1 - q)$. It can be
388 checked that this is a valid transport plan and this matches the upper bound on the optimal cost from
389 Theorem 1.

390 **B.3 Proof of Theorem 1**

391 It would be enough to show that for any $\pi \in \Pi(p^{\otimes k}, q)$, and any $\Omega_0 \subset \Omega$, we have

$$\Pr(Y \in S(X^k)) \leq \sum_{y \in \Omega_0} \min\{q(y), 1 - (1 - p(y))^k\} + \sum_{x^k \in \Omega^k} \min\{\prod_{i=1}^k p(x_i), \sum_{y \in S(x^k) \cap \Omega_0^c} q(y)\}.$$

392

$$\begin{aligned} & \Pr(Y \in S(X^k)) \\ &= \sum_{y \in \Omega} \sum_{x^k \in \Omega^k} \Pr(X^k = x^k, Y = y) \cdot \mathbb{1}\{y \in S(x^k)\} \\ &= \sum_{y \in \Omega_0} \sum_{x^k \in \Omega^k} \Pr(X^k = x^k, Y = y) \cdot \mathbb{1}\{y \in S(x^k)\} \\ &+ \sum_{y \in \Omega_0^c} \sum_{x^k \in \Omega^k} \Pr(X^k = x^k, Y = y) \cdot \mathbb{1}\{y \in S(x^k)\} \\ &= \sum_{y \in \Omega_0} \Pr(y \in S(x^k), Y = y) + \sum_{x^k \in \Omega^k} \sum_{y \in S(x^k) \cap \Omega_0^c} \Pr(X^k = x^k, Y = y) \\ &\leq \sum_{y \in \Omega_0} \min\{\Pr(y \in S(x^k)), q(y)\} + \sum_{x^k \in \Omega^k} \min\{\Pr(X^k = x^k), \sum_{y \in S(x^k) \cap \Omega_0^c} q(y)\} \\ &= \sum_{y \in \Omega_0} \min\{1 - (1 - p(y))^k, q(y)\} + \sum_{x^k \in \Omega^k} \min\{\prod_{i=1}^k p(x_i), \sum_{y \in S(x^k) \cap \Omega_0^c} q(y)\}. \end{aligned}$$

393 **C Proof of Theorem 2 and Theorem 3**

394 **C.1 Proof of Theorem 2**

395 We start by proving the following lemma on γ^* .

396 **Lemma 3.** *Let*

$$f(\gamma) = 1 - (1 - \beta_{p,q}(\gamma))^k - \gamma \beta_{p,q}(\gamma).$$

397 *Then we have Let γ^* be the solution to Eq. (6). Then when $d_{\text{TV}}(p, q) \in (0, 1)$,*

- 398 • $f(\gamma)$ is monotone in γ in $[1, \infty)$;
- 399 • $\gamma^* \in [1, \min\{k, \max_x \frac{q(x)}{p(x)}\}]$.

400 *Proof.* It would enough to prove the followings: (1) $f(\gamma)$ is monotone in γ in $[1, \infty)$; (2) $f(1) \geq 0$;

401 (3) $f(k) \leq 0$; (4) $f(\max_x \frac{q(x)}{p(x)}) \leq 0$.

402 To see (1), since $\beta_{p,q}(\gamma)$ is decreasing in γ , so is $1 - (1 - \beta_{p,q}(\gamma))^k$. Moreover, $\gamma \beta_{p,q}(\gamma) =$
 403 $\sum_x \min\{\gamma p(x), q(x)\}$, which is non-decreasing in γ . Hence we have $1 - (1 - \beta_{p,q}(\gamma))^k - \gamma \beta_{p,q}(\gamma)$
 404 is decreasing.

405 To see (2), note that when $\gamma = 1$, $\beta_{p,q}(\gamma) = 1 - d_{\text{TV}}(p, q)$. Hence we have

$$1 - (1 - \beta_{p,q}(1))^k = 1 - d_{\text{TV}}(p, q)^k \geq 1 - d_{\text{TV}}(p, q).$$

406 When $\gamma = k$, (3) holds since for $x \in [0, 1]$, we have $1 - (1 - x)^k \leq kx$. Moreover, when
 407 $\gamma = \max_x \frac{q(x)}{p(x)} > 1$, we have $\beta_{p,q}(\gamma) = 1/\gamma$ and (4) holds since

$$1 - (1 - \beta_{p,q}(\gamma))^k = 1 - (1 - 1/\gamma)^k < 1 = \gamma \cdot 1/\gamma.$$

408

□

409 Next we prove Theorem 2, we will break the proof into four parts: (1) computation efficiency; (2)
 410 $\pi_\gamma^{\text{K-SEQ}}$ is a valid transport plan; (3) acceptance probability; (4) optimality guarantee of $\pi_\gamma^{\text{K-SEQ}}$.

411 **Computation efficiency.** Note that the lemma immediately implies that γ^* can be computed up to
 412 arbitrary accuracy δ in time $|\Omega| \log(k/\delta)$ using binary search over $[1, k]$.

413 **Valid transport plan.** We next prove that $\pi_\gamma^{\text{k-SEQ}}$ is a valid transport plan when $\gamma \geq \gamma^*$. By
 414 Lemma 3, when $\gamma \geq \gamma^*$, we have $1 - (1 - \beta_{p,q}(\gamma))^k \geq \gamma \beta_{p,q}(\gamma)$. Recall that $p_{\text{acc}} = 1 - (1 - \beta_{p,q}(\gamma))^k$,
 415 and

$$\forall x \in \Omega, p^{\text{res}}(x) = \frac{q(x) - \min \left\{ p(x), \frac{q(x)}{\gamma} \right\} \frac{p_{\text{acc}}}{\beta_{p,q}(\gamma)}}{1 - p_{\text{acc}}}.$$

416 $\forall x \in \Omega$, we have

$$\min \left\{ p(x), \frac{q(x)}{\gamma} \right\} \frac{p_{\text{acc}}}{\beta_{p,q}(\gamma)} \leq \frac{1 - (1 - \beta_{p,q}(\gamma))^k}{\gamma \beta_{p,q}(\gamma)} q(x) \leq q(x),$$

417 this implies $p^{\text{res}}(x) \geq 0$ for all $x \in \Omega$. Moreover,

$$\sum_{x \in \Omega} p^{\text{res}}(x) = \sum_{x \in \Omega} \frac{q(x) - \min \left\{ p(x), \frac{q(x)}{\gamma} \right\} \frac{p_{\text{acc}}}{\beta_{p,q}(\gamma)}}{1 - p_{\text{acc}}} = 1.$$

418 Hence p^{res} is a valid distribution. It remains to show that the marginal of Y is q . We first compute the
 419 probability of the output $Y = x$. Note that probability that $Y = X_1$ is

$$p(X_1) \min \left(1, \frac{q(X_1)}{\gamma p(X_1)} \right) = \min \left(p(X_1), \frac{q(X_1)}{\gamma} \right).$$

420 Hence

$$\Pr(Y = X_1 = x) = \min \left(p(x), \frac{q(x)}{\gamma} \right).$$

421 Therefore,

$$\Pr(Y = X_1) = \sum_x \min \left(p(x), \frac{q(x)}{\gamma} \right) = \beta(\gamma).$$

422 Similarly, probability that

$$\Pr(Y = X_2 = x) = \Pr(Y \neq X_1) \Pr(Y = X_2 | Y \neq X_1) = (1 - \beta_{p,q}(\gamma)) \min \left(p(x), \frac{q(x)}{\gamma} \right).$$

423 Hence,

$$\begin{aligned} & \Pr(Y = x, \text{one of } X^k \text{ is accepted}) \\ &= \sum_{i=0}^{k-1} \Pr(X_1, \dots, X_i \text{ are rejected}, X_{i+1} \text{ is accepted, and } X_{i+1} = x) \\ &= \sum_{i=0}^{k-1} (1 - \beta_{p,q}(\gamma))^i \cdot p(x) \cdot \min \left\{ 1, \frac{q(x)}{p(x)\gamma} \right\} \\ &= \min \left\{ p(x), \frac{q(x)}{\gamma} \right\} \cdot \sum_{i=0}^{k-1} (1 - \beta_{p,q}(\gamma))^i \\ &= \min \left\{ p(x), \frac{q(x)}{\gamma} \right\} \frac{1 - (1 - \beta_{p,q}(\gamma))^k}{\beta_{p,q}(\gamma)} \end{aligned}$$

424 Summing over all symbols x yields

$$\Pr(\text{one of } X^k \text{ is accepted}) = \sum_x \min \left\{ p(x), \frac{q(x)}{\gamma} \right\} \frac{1 - (1 - \beta_{p,q}(\gamma))^k}{\beta_{p,q}(\gamma)} = 1 - (1 - \beta_{p,q}(\gamma))^k.$$

425 Hence we have

$$\begin{aligned}
\Pr(Y = x) &= \Pr(Y = x, \text{one of } X^k \text{ is accepted}) + (1 - \text{one of } X^k \text{ is accepted})p^{\text{res}}(x) \\
&= \min\left\{p(x), \frac{q(x)}{\gamma}\right\} \frac{1 - (1 - \beta_{p,q}(\gamma))^k}{\beta_{p,q}(\gamma)} \\
&\quad + (1 - (1 - \beta_{p,q}(\gamma))^k) \frac{q(x) - \min\left\{p(x), \frac{q(x)}{\gamma}\right\} \frac{1 - (1 - \beta_{p,q}(\gamma))^k}{\beta_{p,q}(\gamma)}}{1 - (1 - \beta_{p,q}(\gamma))^k} \\
&= q(x).
\end{aligned}$$

426 **Acceptance probability.** The acceptance probability holds since

$$\alpha(\pi_{\gamma}^{\text{K-SEQ}}) \geq \Pr(\text{one of } X^k \text{ is accepted}) = 1 - (1 - \beta_{p,q}(\gamma))^k.$$

427 **Optimality guarantee of $\pi_{\gamma^*}^{\text{K-SEQ}}$.** It can be seen that $\beta(\gamma)$ is decreasing in γ , and so is $1 - (1 - \beta_{p,q}(\gamma))^k$. Hence we have

$$\alpha(\pi_{\gamma^*}^{\text{K-SEQ}}) \geq 1 - (1 - \beta_{p,q}(\gamma^*))^k \geq 1 - (1 - \beta_{p,q}(k))^k = c_k(p, q) \cdot \min\{kp(x), q(x)\},$$

429 where

$$c_k(p, q) = \frac{1 - (1 - \beta_{p,q}(k))^k}{k\beta_{p,q}(k)} \in [1 - (1 - 1/k)^k, 1).$$

430 The inclusion holds since $f(x) = \frac{1 - (1-x)^k}{kx}$ is monotonically decreasing when $x \in (0, 1/k]$ and
431 $f(1/k) = 1 - (1 - 1/k)^k$, $\lim_{x \rightarrow 0^+} f(x) = 1$.

432 Moreover, $\forall x \geq 0, kx \geq 1 - (1 - x)^k$. Hence we have

$$\begin{aligned}
\alpha(\pi_{\gamma^*}^{\text{K-SEQ}}) &\geq (1 - (1 - 1/k)^k) \cdot \min\{kp(x), q(x)\} \\
&\geq (1 - (1 - 1/k)^k) \min\{1 - (1 - p(x))^k, q(x)\} \\
&\geq (1 - (1 - 1/k)^k) \alpha_k(p, q),
\end{aligned}$$

433 where the last inequality is due to the upper bound in Theorem 1 with $\Omega_0 = \Omega$.

434 C.2 Proof of Theorem 3

435 We prove the theorem via induction. When $L = 1$, $\Pr(\tau = 1) = 1$, the theorem follows directly
436 since f_π in Algorithm 3 is a valid transport plan. Suppose the theorem holds for $L = \ell \geq 1$, we next
437 prove that it holds for $L = \ell + 1$. Let $\bar{Y}^{\tau'}$ be the output sequence when $L = \ell + 1$ and $\bar{Z}^{\tau'+1:\ell+1}$ be
438 the subsequent samples from \mathcal{M}_b . Note that compared to the case when $L = \ell$, extending the block
439 length of the tree by one only changes the probability of Y^τ when $\tau = L$, i.e., $\forall j < \ell$ and length- j
440 sequence $o^j \in \Omega^j$, we have

$$\Pr(Y^j = o^j, \tau = j) = \Pr(\bar{Y}^j = o^j, \tau' = j)$$

441 For any length- ℓ sequence o^ℓ , let

$$\delta(o^\ell) := \Pr(Y^\ell = o^\ell, \tau = \ell) - \Pr(\bar{Y}^\ell = o^\ell, \tau' = \ell).$$

442 Then by definition, we have

$$\delta(o^\ell) = \sum_{o^{(\ell+1)} \in \Omega} \Pr(Y^{\ell+1} = (o^\ell, o^{(\ell+1)}), \tau = \ell + 1)$$

443 For any length- $(\ell + 1)$ sequence $o^{\ell+1} \in \Omega^{\ell+1}$, we have

$$\begin{aligned}
& \Pr\left((\bar{Y}^j, \bar{Z}^{\tau'+1:\ell+1}) = o^{\ell+1}\right) \tag{9} \\
&= \sum_{j=1}^{\ell-1} \Pr(\bar{Y}^j = o^j, \tau' = j) \mathcal{M}_b(o^{j+1:\ell+1} | x^n, o^j) \\
&\quad + \Pr(\bar{Y}^\ell = o^\ell, \tau' = \ell) \mathcal{M}_b(o(\ell+1) | x^n, o^\ell) + \Pr(\bar{Y}^{\ell+1} = o^{\ell+1}, \tau = \ell + 1) \tag{10} \\
&= \sum_{j=1}^{\ell-1} \Pr(Y^j = o^j, \tau' = j) \mathcal{M}_b(o^{j+1:\ell+1} | x^n, o^j) \\
&\quad + (\Pr(Y^\ell = o^\ell, \tau' = \ell) - \delta(o^\ell)) \mathcal{M}_b(o(\ell+1) | x^n, o^\ell) + \Pr(\bar{Y}^{\ell+1} = o^{\ell+1}, \tau = \ell + 1) \tag{11} \\
&= \mathcal{M}_b(o^{\ell+1} | x^n) - \delta(o^\ell) \mathcal{M}_b(o(\ell+1) | x^n, o^\ell) + \Pr(\bar{Y}^{\ell+1} = o^{\ell+1}, \tau = \ell + 1). \tag{12}
\end{aligned}$$

444 Hence it would enough to show that

$$\delta(o^\ell) \mathcal{M}_b(o(\ell+1) | x^n, o^\ell) = \Pr(\bar{Y}^{\ell+1} = o^{\ell+1}, \tau = \ell + 1) \tag{13}$$

445 Note that the event $\bar{Y}^{\ell+1} = o^{\ell+1}, \tau = \ell + 1$ only happens when o^ℓ are all accepted samples from
446 \mathcal{M}_s in the sampling process and when proceeding, the next obtained token is $o(\ell + 1)$.

447 On the other hand, the $\delta(o^\ell)$ is the probability of the event that the sampling process stops at o^ℓ when
448 $L = \ell$ and proceeds when $L = \ell + 1$, which, by definition of the algorithm, happens if and only if o^ℓ
449 are all accepted samples from \mathcal{M}_s . Moreover, when proceeding, since f_π is a valid transport plan,
450 we have that the next sample is generated from $\mathcal{M}_b(\cdot | x^\ell, o^\ell)$. And hence Eq. (13) holds.

451 This concludes the proof.

452 D Comparisons between OTM- k and κ -SEQ

453 D.1 Examples where the approximate algorithm is optimal

454 In this section, we show that for the example in Figures 1, κ -SEQ achieves the optimal acceptance
455 accuracy. In this case, $p = U(d)$ and $q = U(d/r)$. Recall that the optimal acceptance probability is

$$\alpha_k(U(d), U(d/r)) = 1 - (1 - 1/r)^k.$$

456 For $U(d)$ and $U(d/r)$, we have

$$\beta(\gamma) = \sum_{x \in [d]} \min\{p(x), q(x)/\gamma\} = \frac{1}{\max\{r, \gamma\}}.$$

457 And hence solving $1 - (1 - \beta(\gamma))^k = \gamma \beta(\gamma)$ gives $\gamma^* = r(1 - (1 - 1/r)^k)$. And be Theorem 2, we
458 have

$$\alpha(\pi_{\gamma^*}^{\kappa\text{-SEQ}}) \geq 1 - (1 - \beta(\gamma^*))^k = 1 - (1 - 1/r)^k.$$

459 And the equality holds since this an upper bound for any coupling.

460 D.2 Gap between OTM- k and κ -SEQ

461 To see how OTM- k and κ -SEQ compare in general, we numerically compute the acceptance proba-
462 bility for a pair of compressed conditional distributions. We feed the prompt

463 *“He said he also has asked prosecutors to”*

464 to both large and small models used in Section 8 and obtain the conditional distributions p, q . To
465 make the computation feasible for OTM- k , we take the set of top 10 elements from p, q respectively
466 and set the support S to be the union of the two sets. Then we set p' and q' to be the normalized
467 distribution of p and q over the set S .

468 We then numerically compute the acceptance probability for the optimal transport solution in Section 5
 469 and the approximate solution in Section 6 with different k 's. The result is shown in Fig. 5. When
 470 $k = 1$, the acceptance probability is equal to $1 - d_{TV}(p', q')$ for both solutions. The acceptance
 471 probability increases for both methods as k increases and there exists a gap between the optimal
 472 and approximate solution. We would expect the gap to exist for general conditional distributions
 473 from language models. We leave exploring computationally efficient ways to close this gap as an
 interesting future direction.

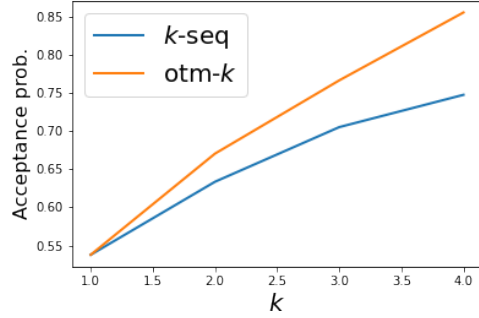


Figure 5: Acceptance probability comparison OTM- k and K-SEQ with compressed conditional distributions.

474

475 E Construct a candidate set by sampling from a prefix-tree

476 As discussed in Section 1, the size of the draft set S is constrained by the number of parallel
 477 computations that can be supported in the hardware. Hence it is important to design the draft set
 478 carefully to allow for a longer sequence of accepted candidate sets. In addition to the *i.i.d.* draft set
 479 selection approach listed in Section 7, we present an algorithm that samples a draft set that forms
 480 the leaves of a prefix tree. Given a draft set size K , the algorithm can be specified by a sequence of
 481 parameter (k_1, k_2, \dots, k_L) satisfying $\prod_{i=1}^L k_i = K$.

482 At a high-level, the algorithm starts with a root node with sequence $x^{1:t}$ and forms a prefix tree of
 483 depth L . At depth $i \in [1 : L - 1]$, each node is expanded by a factor of k_{i+1} and each of its children
 484 will contain a sequence that satisfies: (1) Its prefix agrees with the sequence in the parent node; (2)
 485 The next token is sampled from the conditional probability given the prefix in small model. These
 486 child nodes will be at depth $i + 1$ and the process goes until it hits depth L . We give a detailed
 487 description of the algorithm in Algorithm 5.

Algorithm 5 Draft set selection via prefix-tree.

Input: Input sequence x^t ; expansion factors at each level: (k_1, k_2, \dots, k_L) .

- 1: $S_0 = \{x^t\}$.
 - 2: **for** $i = 0, 1, 2, \dots, L - 1$ **do**
 - 3: $S_{i+1} = \emptyset$.
 - 4: **for all** $\text{seq} \in S_i$ **do**
 - 5: Sample k_{i+1} *i.i.d.* tokens $X_1, X_2, \dots, X_{k_{i+1}}$ from $\mathcal{M}_b(\cdot | \text{seq})$.
 - 6: $S_{i+1} = S_{i+1} \cup \{(\text{seq}, X_i), i = 1, 2, \dots, k_{i+1}\}$.
 - 7: **end for**
 - 8: **end for**
 - 9: **Return** S_L .
-

488 F Additional experiments

489 Similar to Table 1, we report relative latency when parallelizing across the time and batch axes
 490 using the small $6M$ draft model in Table 3. In Table 3, the reported relative latencies are relative to

Table 3: Average latency with parallelization along the time axis and batch axis. We report average latency with standard deviation over 1,000 runs using a 6M transformer relative to the 97M transformer with length = 4 and batch = 1 on GPU.

Relative latency	batch = 1	batch = 2	batch = 4	batch = 8
length = 4	0.18 ± 0.02	0.19 ± 0.04	0.18 ± 0.09	0.20 ± 0.13
length = 8	0.17 ± 0.04	0.19 ± 0.05	0.16 ± 0.02	0.18 ± 0.04

Table 4: Experimental results on the LM1B dataset with varying draft model sizes and the 97M transformer as the large model. All results are over 1000 test prompts averaged over three different random seeds and sampling temperature of 1.0 for both the draft and large models.

Draft model	Algorithm	K	L	Number of decoded tokens per serial call	
2M Transformer	Baseline	-	-	1.00	
	Speculative	1	4	1.86 ± 0.02	
	SpecTr	2	4	2.07 ± 0.01	
	SpecTr	4	4	2.32 ± 0.00	
	SpecTr	8	4	2.56 ± 0.01	
	Speculative	1	8	1.91 ± 0.01	
	SpecTr	2	8	2.15 ± 0.01	
	SpecTr	4	8	2.41 ± 0.00	
	SpecTr	8	8	2.68 ± 0.01	
	6M Transformer	Baseline	-	-	1.00
		Speculative	1	4	2.21 ± 0.01
		SpecTr	2	4	2.43 ± 0.01
SpecTr		4	4	2.74 ± 0.01	
SpecTr		8	4	2.99 ± 0.02	
Speculative		1	8	2.33 ± 0.01	
SpecTr		2	8	2.61 ± 0.02	
SpecTr		4	8	2.96 ± 0.03	
SpecTr		8	8	3.27 ± 0.02	
20M Transformer		Baseline	-	-	1.00
		Speculative	1	4	2.71 ± 0.01
		SpecTr	2	4	2.96 ± 0.00
	SpecTr	4	4	3.28 ± 0.02	
	SpecTr	8	4	3.49 ± 0.03	
	Speculative	1	8	3.12 ± 0.02	
	SpecTr	2	8	3.48 ± 0.04	
	SpecTr	4	8	3.85 ± 0.05	
	SpecTr	8	8	4.15 ± 0.04	

491 the large 97M model to get a sense of the relative cost of sampling multiple drafts with the small
492 model compared to the large model. We also include results for varying draft model sizes with
493 the same 97M large model for LM1B in Table 4. These additional draft models were produced by
494 either halving (2M) or doubling (20M) the original 6M draft model’s number of layers, embedding
495 dimension, MLP dimension, and number of attention heads. As expected, the larger draft models
496 improve all speculative methods’ number of decoded tokens per large model serial call with SpecTr
497 maintaining the best performance across all draft model sizes.