

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 GREEDY MULTI-PATH BLOCK VERIFICATION FOR FASTER DECODING IN SPECULATIVE SAMPLING

Anonymous authors

Paper under double-blind review

ABSTRACT

The goal of L -step speculative decoding is to accelerate autoregressive decoding of a target model by using a cheaper draft model to generate a candidate path of L tokens. Based on a verification algorithm involving target and draft model probabilities, a prefix of the candidate sequence is accepted, and an additional correction token is sampled from a residual distribution to ensure that the final output adheres to the target distribution. While standard speculative decoding uses a verification algorithm which is independent at each token on the path, a recent extension called block verification uses a joint condition involving all sampled on-path probabilities. Block verification (BV) was shown to be optimal over all verification algorithms which use only on-path probabilities, improving on standard speculative decoding. In this work, we first show that block verification is optimal even over verification algorithms that use off-path probabilities, by constructing an information-agnostic linear program (LP). Further, we can extend our LP to the setting where the draft model samples multiple candidate paths, and use it to construct a natural class of multi-path block verification generalizations. While computing the optimal algorithm in this class is not tractable, by considering a stricter class of greedy algorithms, we can formulate an efficient method called greedy multi-path block verification (GBV). Empirically, GBV can improve block efficiency by over 30% and reduce decoding walltimes by over 15% relative to BV.

1 INTRODUCTION

Large language models (LLMs) achieve strong results across code, language, and reasoning (Zhu et al., 2024; Kasneci et al., 2023; Thirunavukarasu et al., 2023). Most LLM families, such as Qwen (Bai et al., 2023), GPT (Radford et al., 2018; 2019; Brown et al., 2020; OpenAI, 2023), and Llama (Touvron et al., 2023b;a) employ autoregressive decoding. When inference is performed with these transformer-based architectures on GPUs, end-to-end latency is dominated by memory bandwidth rather than compute (Fu et al., 2024).

Speculative sampling (Chen et al., 2023; Leviathan et al., 2023) aims to reduce such punitive costs when sampling from a large *target model*. This procedure autoregressively decodes a candidate sequence from a cheaper *draft model*, performs a forward pass over the candidate sequence with the target model to obtain target distribution values, and then probabilistically alters the sequence through a *verification algorithm* to ensure that the resulting output matches the target model distribution. Because inference is bandwidth-bound, the target forward pass over the candidate sequence presents negligible overhead over a standard forward pass. Thus, speculative sampling can decode many tokens in a single target model call and speed up inference without affecting downstream performance.

In **standard speculative sampling**, the draft model proposes a length L draft block, and the verification algorithm independently accepts or rejects each token based on target and draft model distributions values along the draft block. The longest prefix with no rejections is chosen, and an additional token is sampled from a residual distribution, ensuring that at least one token is always generated. While this procedure significantly improves decoding efficiency, it suffers when there are low acceptance rates at early tokens. If the first token is almost always rejected due to a poor draft suggestion, then even if subsequent tokens are always accepted, there is no speedup.

To overcome this early-token bottleneck, recent work on **block verification (BV)** (Sun et al., 2024b) replaces independent token-wise verification with independent prefix-wise verification, and selects

054 the longest accepted prefix. An additional token is still sampled, but the residual distribution is
 055 altered to ensure the output still matches the target distribution. This ensures that more tokens can be
 056 accepted when earlier tokens are likely to be rejected. Block verification is empirically 5-8% faster
 057 than standard speculative sampling. In fact, it is provably optimal over all verification algorithms that
 058 only use target and draft probability distributions along the drafted candidate sequence, including
 059 standard speculative sampling.

060 In this paper, we frame the optimal verification algorithm as a solution to a linear program (LP), and
 061 extend the LP to the setting of *multiple* i.i.d. generated draft sequences. Our contributions are:
 062

- 063 • **Single-Path Information-Agnostic LP.** We develop an LP that encodes feasible speedups
 064 from verification algorithms that must still return a candidate sequence prefix and an
 065 additional token, but are now given access to the *full joint* target and draft distributions, not
 066 just on-path values. We show that block verification is still optimal in this setting. That is,
 067 the prefix output requirement is the bottleneck to improving optimal decoding efficiency,
 068 rather than off-path information access.
- 069 • **Multi-Path Information-Agnostic LP.** We extend the single-path LP to the setting where the
 070 draft model i.i.d. samples $K > 1$ candidate sequences, considering verification algorithms
 071 which return a prefix of one of the paths and an additional token. Unlike the $K = 1$
 072 setting, we find that knowledge of off-path target and draft distributions can improve
 073 decoding efficiency. We show that the optimal verification algorithm in this setting is to
 074 probabilistically select one of the K paths, and then run block verification on that single
 075 path with a skewed draft distribution.
- 076 • **Greedy Approximation Schemes.** The optimal solution to the multi-path LP involves a
 077 complex nonlinear optimization problem over the joint target and draft distributions, and is
 078 infeasible to solve directly. Thus, we explore a class of greedy approximations: algorithms
 079 which globally rank all possible candidate sequences, select the highest-ranking of the K
 080 paths, and run block verification on that path. We show that these can be viewed as outputs
 081 of a greedy polymatroid algorithm in the multi-path LP.
- 082 • **Greedy Multi-Path Block Verification (GBV).** Many greedy approximation schemes
 083 require the full joint target and draft distributions. However, when the global ranking is
 084 tree-based, only on-path values are required. We devise a simple tree-based rule that results
 085 in wall-clock speedups of over 15% relative to BV, and provide theoretical justification for
 086 this rule.

087 2 BACKGROUND

089 We first review standard speculative sampling and its extensions. We start with the case where only
 090 one draft block is generated, and then cover multi-path extensions.
 091

092 2.1 SINGLE-PATH

094 Denote the target model by M_p , and the draft model by M_q , which are assumed to share a common
 095 vocabulary \mathcal{V} . We use the sequence notation $a_{1:k} = (a_1, \dots, a_k)$ for tokens $a_1, \dots, a_k \in \mathcal{V}$ and
 096 the inclusive slicing notation $a_{i:j} = (a_i, \dots, a_j)$, with the empty sequence convention $a_{i:j} = \emptyset$ for
 097 $j < i$. Given a context c , the target and draft models induce next-token distribution $p(\cdot|c)$ and $q(\cdot|c)$,
 098 respectively. In autoregressive sampling, these further induce next- k -token distributions through
 099 conditional factorization, which we denote:

$$100 \quad q_k(a_{1:k}|c) = \prod_{i=1}^k q(a_i|c, a_{1:i-1}), \quad p_k(a_{1:k}|c) = \prod_{i=1}^k p(a_i|c, a_{1:i-1}). \quad (1)$$

103 In L -step speculative sampling, we autoregressively sample a draft block of L tokens from the draft
 104 distribution q_L and call the target model M_p on this block *once*, to obtain target and draft next-token
 105 probabilities along the block. Through a randomized rule called the **verification algorithm**, we
 106 accept the first $\tau \in \{0, \dots, L\}$ of these L tokens and sample an additional correction token, such
 107 that this output matches the true target distribution. The **block efficiency** $\mathbb{E}[\tau + 1]$ is the average
 number of decoded tokens per M_p call; in the special case where the target and draft distributions are

108 identical, it achieves its maximum value of $L + 1$. When a cheap draft model is used and L is not too
 109 large, block efficiency is an accurate indicator of walltime speedup.
 110

111 **Speculative sampling.** The original schemes of Chen et al. (2023); Leviathan et al. (2023) first
 112 generate a draft block $a_{1:L} \sim q_L(\cdot|\mathbf{c})$. They *independently* accept each token a_i with probability
 113 $\min(1, p(a_i|a_{1:i-1})/q(a_i|a_{1:i-1}))$, and τ is the maximum value such that $a_{1:\tau}$ consists of only
 114 accepted tokens. Essentially, they choose the longest prefix of acceptances. Finally, they sample an
 115 additional *correction token* $y \sim p_{\text{res}}^{\text{token}}(\cdot|\mathbf{c}, a_{1:i})$ from a residual distribution if $\tau < L$:

$$p_{\text{res}}^{\text{token}}(\cdot|\mathbf{c}, a_{1:i}) \propto \max\{p(\cdot|\mathbf{c}, a_{1:i}) - q(\cdot|\mathbf{c}, a_{1:i}), 0\}. \quad (2)$$

116 If $\tau = L$, they just sample $y \sim p(\cdot|\mathbf{c}, a_{1:L})$ directly from the target. This ensures that the final output
 117 $(a_{1:\tau}, y)$ matches the target distribution.
 118

119 **Draft extensions.** Some recent works have improved upon block efficiency in speculative sampling
 120 by altering the drafting phase, through retrieval or cascading (He et al., 2023; Chen et al., 2024),
 121 hierarchical drafting (Sun et al., 2024a), distillation (Zhou et al., 2023; Liu et al., 2023), layer skipping
 122 (Zhang et al., 2023; Elhoushi et al., 2024), or multi-token heads (Gloeckle et al., 2024; Samragh et al.,
 123 2025)). In this paper, we instead focus on methods which improve block efficiency by altering the
 124 verification algorithm, and these can be integrated with any of the above works for further gains.
 125

126 **Tree verification.** Monte Carlo tree verification from Hu & Huang (2024) provably improves the
 127 block efficiency of standard speculative sampling. While the underlying tree structure might suggest
 128 this is a multi-path method, we emphasize that it is a single-path algorithm, as the end of Section 5 in
 129 Hu & Huang (2024) specifies a linear chain of draft tokens. This means it is provably worse than
 130 block verification: see the third paragraph of Section 7 in Sun et al. (2024b).
 131

132 **Block verification (BV).** To improve block efficiency in the verification stage, Sun et al. (2024b)
 133 relax the token-independence assumption for acceptance in standard speculative sampling. First, they
 134 sample the draft block $a_{1:L} \sim q_L(\cdot|\mathbf{c})$, and recursively define on-path weights w :

$$w(\emptyset|\mathbf{c}) = 1, \quad w(a_{1:i}|\mathbf{c}) = \min\left\{1, \frac{w(a_{1:i-1}|\mathbf{c})p(a_i|a_{1:i-1})}{q(a_i|a_{1:i-1})}\right\}. \quad (3)$$

135 Now, they independently accept each *prefix* $a_{1:i}$ with probability

$$h^{\text{block}}(a_{1:i}|\mathbf{c}) = \frac{\sum_{x \in \mathcal{V}} \max\{p(x|\mathbf{c}, a_{1:i}) - q(x|\mathbf{c}, a_{1:i}), 0\}}{1 - w(a_{1:i}|\mathbf{c}) + \sum_{x \in \mathcal{V}} \max\{p(x|\mathbf{c}, a_{1:i}) - q(x|\mathbf{c}, a_{1:i}), 0\}} \quad (4)$$

136 for $i < L$, and $a_{1:L}$ with probability $w(a_{1:L}|\mathbf{c})$. They select the longest accepted prefix, of length τ ,
 137 and sample a correction token $y \sim p_{\text{res}}^{\text{block}}(\cdot|\mathbf{c}, a_{1:i})$, a weighted version of $p_{\text{res}}^{\text{token}}$, if $\tau < L$:

$$p_{\text{res}}^{\text{block}}(\cdot|\mathbf{c}, a_{1:i}) \propto \max\{w(a_{1:i}|\mathbf{c})p(\cdot|\mathbf{c}, a_{1:i}) - q(\cdot|\mathbf{c}, a_{1:i}), 0\}. \quad (5)$$

138 If $\tau = L$, they sample $y \sim p(\cdot|\mathbf{c}, a_{1:L})$. Like in speculative sampling, the output $(a_{1:\tau}, y)$ follows
 139 the target distribution. Sun et al. (2024b) prove that block verification achieves the highest block
 140 efficiency among any verification algorithm that only takes in on-path distribution values.
 141

142 2.2 MULTI-PATH

143 There are also recent lines of work that extend L -step speculative sampling to the *multi-path* setting
 144 (Sun et al., 2023; Spector & Re, 2023; Cai et al., 2024; Li et al., 2024). In this setting, $K > 1$
 145 length- L draft blocks are sampled from q_L . The target model M_p is again called *once* on all draft
 146 blocks in parallel, to obtain target and draft next-token probabilities along all K paths. Then, using
 147 a **multi-path verification algorithm**, a prefix of one of the blocks is accepted, and an additional
 148 correction token from a residual distribution is sampled in order to match the target distribution.
 149 Block efficiency is defined in the same way as previously. For moderate K and L , leveraging GPU
 150 parallelization, there is generally little overhead in performing the *batched* forward pass across K
 151 sequences, relative to a forward pass on one sequence (Agrawal et al., 2024; Dao et al., 2022).
 152

153 We use the concept of a **draft tree** to enumerate all distinct prefixes among the K sampled paths. We
 154 will compare our methods against two categories of multi-path verification algorithms which select a
 155 node from the draft tree, i.e. a prefix of one of the blocks, in different ways.
 156

162 **OTLP-based algorithms.** These perform a top-down traversal of the draft tree. At each step, they
 163 compute a feasible solution of an optimal transport linear program (OTLP) to determine the next
 164 token and progress to a child node, terminating when they land off the draft tree. SpecTr (Sun et al.,
 165 2023), NSS (Miao et al., 2024), and SpecInfer (Miao et al., 2024) are key examples, and recent
 166 theoretical OTLP solvers (Khisti et al., 2025; Hu et al., 2025) further improve block efficiency.
 167

168 **Traversal verification.** To the best of our knowledge, traversal verification in Weng et al. (2025) is
 169 the only multi-path algorithm that traverses the draft tree from the bottom-up. Like our algorithm in
 170 Section 5, this reduces to BV in the single-path $K = 1$ case. In our experiments, traversal verification
 171 and our methods outperform OTLP-based methods, but are each effective in different (p, q) regimes.
 172

173 3 SINGLE-PATH INFORMATION-AGNOSTIC LP

175 We first formally define the class of single-path verification algorithms. We denote $\mathcal{V}^{\leq k} = \mathcal{V}^0 \cup$
 176 $\mathcal{V}^1 \cup \dots \cup \mathcal{V}^k$, where \mathcal{V}^k is the set of length- k sequences of tokens in \mathcal{V} , and thus $\mathcal{V}^{\leq k}$ is the set of
 177 sequences of length $\leq k$. We use L to denote the draft block length. For the remainder of the paper,
 178 we use notation from Section 2, omitting the context c and subscripts on p, q when they are clear.

179 **Definition 3.1** (Single-path draft verification algorithm). A **single-path verification algorithm** Φ
 180 takes in a sampled draft block $X_{1:L} \sim q_L(\cdot)$, and the full L -step target and draft distributions p_L and
 181 q_L , and returns a nonempty sequence in $\mathcal{V}^{\leq L+1}$. It is **valid** under:

- 182 • **Prefix-matching:** the algorithm returns a (possibly empty) prefix $X_{1:\tau}$ of the draft block
 183 followed by an additional correction token Y , for some $\tau \in \{0, \dots, L\}$.
 184
- 185 • **Target-matching:** for any context c , any draft and target models M_p and M_q , and any
 186 block length L , when we generate $L - \tau$ additional tokens $Z \sim p_{L-\tau}(\cdot | X_{1:\tau}, Y)$, we have
 187 $(X_{1:\tau}, Y, Z) \sim_p p_{L+1}(\cdot)$, where \sim_p denotes equality of distributions.

188 If we also require that Φ can only take in the *on-path* target and draft distributions $p(\cdot | X_{1:i})$, $q(\cdot | X_{1:i})$
 189 for $i \in \{0, \dots, L\}$, then it is called **information-restricted**.
 190

191 Target-matching means that sampling from M_p autoregressively after verification, until a total of $L + 1$
 192 tokens are generated, gives the same result as just sampling $L + 1$ tokens from M_p autoregressively
 193 without verification. This is equivalent to guaranteeing that running the verification algorithm and
 194 appending its output to the context iteratively maintains the target distribution, even as τ varies: see
 195 Lemma 2 in Appendix B of Sun et al. (2024b) for a formal proof.

196 Prefix-matching forces the verification output to only deviate from the draft block at its last token.
 197 This is a strict modeling requirement, not a technical convenience. To bypass prefix-matching while
 198 maintaining the target distribution, one must access the next-token target distribution $p(\cdot | c)$ for
 199 contexts c outside the draft tree. This is not possible without performing another M_p forward pass,
 200 which would defeat the goal of speculative decoding: generating multiple tokens in one target call.

201 Importantly, our definition is *less strict* than in Sun et al. (2024b), because they only consider the class
 202 of valid single-path *information-restricted* verification algorithms, and prove that block verification is
 203 optimal in this class. Surprisingly, we find that block verification is optimal in the class of *all* valid
 204 single-path verification algorithms. To prove this, we first define node budgets, which represent how
 205 much mass a verification algorithm has allocated along a path relative to the target distribution.

206 **Definition 3.2.** For any single-path verification algorithm Φ , we define **node budgets**

$$207 \quad D_{\Phi}(a_{1:i}) = 1 - \sum_{j=1}^i \frac{\mathbb{P}(X_{1:\tau} = a_{1:j-1}, Y = a_j)}{p(a_{1:j})} \quad \forall a_{1:i} \in \mathcal{V}^{\leq L+1}. \quad (6)$$

210 This can be represented by a function $D_{\Phi} : \mathcal{V}^{\leq L+1} \rightarrow \mathbb{R}$.
 211

212 Using node budget variables, we define the **single-path information-agnostic LP** in Theorem 3.3,
 213 which describes what values of D_{Φ} a valid single-path verification algorithm Φ can induce. Here,
 214 the chain of D_{Φ} inequalities along paths encode the target-matching constraint, and the pointwise
 215 upper bounds on $D_{\Phi} p$ encode prefix-matching. We defer the proof to Appendix B, as the special case
 $K = 1$ of the multi-path LP in Theorem 4.3, which we prove in Appendix A.

216 **Theorem 3.3.** *A single-path verification algorithm Φ is valid if and only if the node budgets*
 217 $D_\Phi : \mathcal{V}^{\leq L+1} \rightarrow \mathbb{R}$ *are feasible in the following LP:*

$$219 \quad \max \sum_{a_{1:i} \in \mathcal{V}^{\leq L}} D_\Phi(a_{1:i}) p(a_{1:i}) \quad (7)$$

$$221 \quad \text{s.t. } 1 = D_\Phi(a_{1:0}) \geq D_\Phi(a_{1:1}) \geq \dots \geq D_\Phi(a_{1:L}) \geq D_\Phi(a_{1:L+1}) = 0 \quad \forall a_{1:L+1} \in \mathcal{V}^{L+1}, \quad (8)$$

$$223 \quad D_\Phi(a_{1:i}) p(a_{1:i}) \leq q(a_{1:i}) \quad \forall a_{1:i} \in \mathcal{V}^{\leq L}. \quad (9)$$

224 *Furthermore, the objective value is precisely the block efficiency $\mathbb{E}[\tau + 1]$ for Φ .*

226 Because the feasibility conditions in the single-path LP only involve pointwise and pathwise inequalities,
 227 it is not hard to compute the optimal objective value and node budgets for a valid single-path
 228 verification algorithm, by using greedy allocation along paths. These node budgets match those
 229 derived from block verification, and thus block verification is the optimal valid single-path algorithm.

230 **Theorem 3.4.** *Block verification, which is a valid single-path verification algorithm Φ_{BV} , achieves*
 231 *the highest block efficiency of any valid single-path verification algorithm.*

233 See [Appendix C](#) for a proof of [Theorem 3.4](#). These results show that prefix-matching is the key
 234 barrier to improving acceptance: even with access to the full joint target and draft distributions, prefix-
 235 matching (pointwise upper bounds on $D_\Phi p$ terms) prevent us from getting better block efficiency (the
 236 sum of $D_\Phi p$ terms) than block verification. This motivates our exploration of verification algorithms
 237 which draft *multiple* candidate paths in the next section. In this setting, the valid prefix output space
 238 grows, thereby loosening the prefix-matching requirement and improving block efficiency.

4 MULTI-PATH INFORMATION AGNOSTIC LP

242 In this section, we show that off-path distribution information can improve block efficiency when
 243 $K > 1$ draft paths are generated, even with prefix-matching and target-matching requirements. We
 244 first define the class of multi-path verification algorithms, and extend [Theorem 3.3](#) to these algorithms.

245 **Definition 4.1** (K -path draft verification algorithm). A K -path verification algorithm Φ takes in K
 246 i.i.d. sampled draft blocks $X_{1:L}^{(1)}, \dots, X_{1:L}^{(K)} \sim q_L(\cdot | c)$, and the full L -step distributions p_L and q_L ,
 247 and returns a nonempty sequence in $\mathcal{V}^{\leq L+1}$. It is **valid** under:

- 249 • **Prefix-matching:** the algorithm returns a (possibly empty) prefix $X_{1:\tau}$ of *some* draft block
 250 followed by an additional correction token Y , for some $\tau \in \{0, \dots, L\}$.
- 251 • **Target-matching:** this is the same as in [Definition 3.1](#).

253 If Φ can only take in the *on-path* target and draft distributions $p(\cdot | X_{1:i}^{(k)})$, $q(\cdot | X_{1:i}^{(k)})$ for $k \in [K]$ and
 254 $i \in \{0, \dots, L\}$, $k \in [K]$, then it is called **information-restricted**.

256 To the best of our knowledge, all existing valid multi-path verification algorithms (see [Section 2.2](#))
 257 are information-restricted. We are the first to consider theoretical efficiency limits in the absence of
 258 information-restriction, and explicitly encode prefix-matching and target-matching into an LP.

259 To extend the single-path LP from [Theorem 3.3](#) to this setting, we first define node budgets D_Φ
 260 in the same way as [Definition 3.2](#). We also require the notion of an *antichain*, which is a set of
 261 variable-length paths in $\mathcal{V}^{\leq L}$ where no path is a prefix of another. Antichains are useful because the
 262 mass of an autoregressive distribution over an antichain can be computed by summing its probabilities
 263 at all antichain elements, due to the disjointness of autoregressively sampling antichain elements.

264 **Definition 4.2.** An **antichain** of $\mathcal{V}^{\leq L}$ is a subset of $\mathcal{V}^{\leq L}$ where no sequence in the antichain is a
 265 prefix of another. We denote the set of all antichains in $\mathcal{V}^{\leq L}$ by $\mathcal{A}(\mathcal{V}^{\leq L})$.

267 Now, to form the **multi-path information-agnostic LP** in [Theorem 4.3](#), we replace the pointwise
 268 upper bounds on $D_\Phi p$ by subset-sum upper bounds over antichains. As mentioned at the end of
 269 [Section 3](#), this corresponds to altering the prefix output space. The proof uses Hall-type feasibility
 constraints, and is fairly involved. Due to space constraints, we defer the proof to [Appendix A](#).

270 **Theorem 4.3.** A K -path verification algorithm Φ is valid if and only if the node budgets $D_\Phi : \mathcal{V}^{\leq L+1} \rightarrow \mathbb{R}$ are feasible in the following LP:

$$273 \quad \max \sum_{a_{1:i} \in \mathcal{V}^{\leq L}} D_\Phi(a_{1:i}) p(a_{1:i}) \quad (10)$$

$$275 \quad \text{s.t. } 1 = D_\Phi(a_{1:0}) \geq D_\Phi(a_{1:1}) \geq \dots \geq D_\Phi(a_{1:L}) \geq D_\Phi(a_{1:L+1}) = 0 \quad \forall a_{1:L+1} \in \mathcal{V}^{L+1}, \quad (11)$$

$$277 \quad \sum_{a_{1:i} \in T} D_\Phi(a_{1:i}) p(a_{1:i}) \leq 1 - \left(1 - \sum_{a_{1:i} \in T} q(a_{1:i})\right)^K \quad \forall T \in \mathcal{A}(\mathcal{V}^{\leq L}). \quad (12)$$

280 Furthermore, the objective value is precisely the block efficiency $\mathbb{E}[\tau + 1]$ for Φ .

282 This reduces to the single-path LP in the case $K = 1$. However, for $K > 1$, there is an important
283 distinction: the upper bound on sums of $D_\Phi p$ now represent a *submodular* function of T , unlike
284 the *additive* (*modular*) pointwise upper bounds on $D_\Phi p$ in [Theorem 3.3](#). We now explain how this
285 submodularity naturally arises from a **path selection rule** in valid K -path verification algorithms. In
286 [Lemma 4.4](#) (proof in [Appendix D](#)), we first prove a canonical decomposition of such algorithms into
287 randomized selection of one of the drafted K paths, followed by valid single-path verification.

288 **Lemma 4.4.** A valid K -path verification algorithm has the following equivalent definition. First,
289 given K paths sampled i.i.d. from q , randomly select one through a path selection rule Γ , which can
290 depend on off-path probability values. Say this path follows the distribution q^Γ over \mathcal{V}^L . Then, run a
291 valid single-path verification algorithm on this path with target values p and draft values q^Γ .

292 We call q^Γ the **skewed draft distribution** induced by Γ . Because this is a distribution over \mathcal{V}^L , just
293 as for p_L, q_L , we can naturally extend it to an induced distribution over each \mathcal{V}^i for $0 \leq i \leq L$:

$$295 \quad q^\Gamma(a_{1:i}) = \sum_{a_{i+1:L} \in \mathcal{V}^{L-i}} q^\Gamma(a_{1:L}) \quad \forall a_{1:i} \in \mathcal{V}^{\leq L}. \quad (13)$$

297 This further induces an autoregressive distribution in the same way as p, q :

$$299 \quad q^\Gamma(a_{k+1:i} | a_{1:k}) = \frac{q^\Gamma(a_{1:i})}{q^\Gamma(a_{1:k})} \quad \forall a_{1:i} \in \mathcal{V}^{\leq L}. \quad (14)$$

302 Not all distributions can be realized as the skewed draft distribution of a path selection rule, given a
303 fixed draft q and path count K . In fact, as in [Lemma 4.5](#), the realizable distributions are those which
304 satisfy a submodular inequality like that in [Theorem 4.3](#). For a proof, see [Appendix E](#).

305 **Lemma 4.5.** Fix K , a draft distribution q , and a distribution r . Then $r = q^\Gamma$ for some randomized
306 path selection rule Γ over K drafts sampled i.i.d. from q if and only if

$$307 \quad \sum_{a_{1:i} \in T} r(a_{1:L}) \leq 1 - \left(1 - \sum_{a_{1:i} \in T} q(a_{1:i})\right)^K \quad \forall T \in \mathcal{A}(\mathcal{V}^{\leq L}). \quad (15)$$

312 This shows that [Lemma 4.4](#) corresponds to **linearizing** the submodular constraints in [Theorem 4.3](#).
313 Indeed, consider replacing the *submodular* upper bounds in the multi-path LP with the *additive*
314 (*modular*) upper bounds $\sum_{a_{1:i} \in T} q^\Gamma(a_{1:i})$. By the upper bound in [Lemma 4.5](#), any solution to this
315 **linearized LP** will satisfy the original multi-path LP. In fact, this linearized LP is precisely the
316 single-path LP¹ where the draft distribution is the skewed draft q^Γ rather than the original q . This
317 corresponds to running valid single-path verification with target p and draft q^Γ , as in [Lemma 4.4](#).

318 For a given q^Γ , using the result of [Theorem 3.4](#), the optimal valid single-path verification algorithm
319 to run here is precisely block verification. This leads to an explicit description of the optimal valid
320 multi-path verification algorithm in [Theorem 4.6](#). We prove this in [Appendix F](#).

321 **Theorem 4.6.** The optimal valid multi-path verification algorithm randomly chooses one of the K
322 i.i.d. blocks with Γ and then runs single-path block verification on that path with target values p

323 ¹See [Appendix B](#) for an explanation of how the sums reduce to pointwise inequalities like in [Theorem 3.3](#).

324 and draft values q^Γ , for some path selection rule Γ . Furthermore, the optimal choice of Γ can be
 325 determined by solving the following optimization problem:
 326

$$327 \max_{a_{1:i} \in \mathcal{V}^{\leq L}} \sum_{0 \leq k \leq i} p(a_{k+1:i} | a_{1:k}) q^\Gamma(a_{1:k}) \quad (16)$$

$$329 \text{s.t. } \sum_{a_{1:L} \in T} q^\Gamma(a_{1:L}) \leq 1 - \left(1 - \sum_{a_{1:L} \in T} q(a_{1:L})\right)^K \quad \forall T \subseteq \mathcal{V}^L. \quad (17)$$

333 This optimization problem is intractable to solve exactly, since it requires knowledge of the full
 334 joint and target distributions. Furthermore, even if the optimal q^Γ can be computed, there is no
 335 guarantee that we can compute an efficient path selection rule Γ which induces this skewed draft
 336 distribution. Therefore, we turn to approximation-based schemes. As [Lemma 4.7](#) shows, we can
 337 obtain a near-optimal solution by making all ratios q^Γ/p close to one. For a proof, see [Appendix G](#).

338 **Lemma 4.7.** *For a fixed q^Γ , the objective value in [Theorem 4.6](#) is lower bounded by*

$$340 (L+1) \cdot \min_{a_{1:i} \in \mathcal{V}^{\leq L}} \frac{q^\Gamma(a_{1:i})}{p(a_{1:i})} \quad (18)$$

343 This shows that when $q^\Gamma = p$ is feasible in [Theorem 4.6](#), we can get an optimal block efficiency of
 344 $L+1$, i.e. we always accept a full length- L path and an additional token. However, even determining
 345 a rule Γ that achieves q^Γ near p is difficult. Thus, in the next section, we restrict ourselves to a
 346 more tractable class of *greedy* algorithms, where one can both explicitly compute the randomized
 347 rule Γ and the resulting distribution values q^Γ . This explicit characterization is necessary to apply
 348 [Lemma 4.4](#), because block verification requires exact draft probability values q^Γ along the selected
 349 path, and it is impossible to exactly sample from q^Γ without an explicit randomized rule Γ .
 350

351 5 GREEDY MULTI-PATH BLOCK VERIFICATION

353 We now design an explicit path selection rule Γ where the skewed draft distribution q^Γ is heuris-
 354 tically close to p and simple to compute, using only on-path probability values. This leads to our
 355 main algorithm: **greedy multi-path block verification (GBV)**. We first define greedy verification
 356 algorithms.

358 **Definition 5.1.** A **greedy multi-path valid verification algorithm** forms a global ranking on paths
 359 in \mathcal{V}^L , and always selects the highest-ranked drafted path in the path selection rule Γ ([Lemma 4.4](#)).

360 In the context of the multi-path LP, each greedy multi-path valid verification algorithms can be viewed
 361 as the output of a greedy polymatroid algorithm ([Schrijver et al., 2003](#)) on the submodular feasibility
 362 constraints in [Theorem 4.3](#), with the algorithm-selected order being the same as the global ranking.
 363 Due to space constraints, we expand on this connection in detail in [Appendix H](#).

364 If the global ordering of paths is $\mathcal{V}^L = \{P_1, \dots, P_M\}$, then one selects P_i from Γ if and only if all
 365 K drafted paths lie in $\{P_1, \dots, P_i\}$, but do not all lie in $\{P_1, \dots, P_{i-1}\}$. Thus, we can explicitly write:

$$367 q^\Gamma(P_i) = \left(\sum_{j=1}^i q(P_j) \right)^K - \left(\sum_{j=1}^{i-1} q(P_j) \right)^K. \quad (19)$$

371 Again, computing this for arbitrary global orderings may require off-path probability values. However,
 372 when the paths are ordered through a tree-based rule, only on-path values are needed. Tree-based
 373 rules create a global ranking by generating local orderings at all prefixes in $\mathcal{V}^{\leq L}$, and then combine
 374 them to induce a lexicographic ordering over all length- L paths.

375 **Definition 5.2.** A global ranking of paths in \mathcal{V}^L is **tree-based** if the ranking can be obtained as
 376 follows. Define an injective function $\pi_{a_{1:i}} : \mathcal{V} \rightarrow \mathbb{R}$ at each $a_{1:i} \in \mathcal{V}^{\leq L}$, only using the values
 377 $p(\cdot | a_{1:i}), q(\cdot | a_{1:i})$. Then, assign to each path $a_{1:L} \in \mathcal{V}^L$ the L -tuple $O(a_{1:L}) = (\pi_{a_{1:i-1}}(a_i))_{i=1}^L$.
 378 Finally, rank the paths $a_{1:L} \in \mathcal{V}^L$ in increasing lexicographic order of $O(a_{1:L})$.

378

Algorithm 1: Greedy Multi-Path Block Verification380 **Input:** Draft blocks $X_{1:L}^{(1)}, \dots, X_{1:L}^{(K)}$, target and draft probabilities $p(\cdot|X_{1:i}^{(k)}), q(\cdot|X_{1:i}^{(k)})$ 381 **1** **for** $k = 1, \dots, K$ **do**382 **2** Define k th path ranking O_k by $\left[p(X_i^{(k)}|X_{1:i-1}^{(k)})/q(X_i^{(k)}|X_{1:i-1}^{(k)}) \right]_{i=1}^L$ 384 **3** Select k_0 with maximal ranking O_{k_0} by lexicographic ordering385 **4** **for** $i = 1, \dots, L$ **do**386 **5** Compute $q^{\Gamma_g}(\cdot|X_{1:i}^{(k_0)})$ from $p(\cdot|X_{1:i}^{(k_0)}), q(\cdot|X_{1:i}^{(k_0)})$ using [Equation \(23\)](#)388 **6** Run block verification on $X_{1:L}^{(k_0)}$ with target and draft probabilities $p(\cdot|X_{1:i}^{(k_0)}), q^{\Gamma_g}(\cdot|X_{1:i}^{(k_0)})$

389

390

391 Now, the question remains of how to choose the local orderings $\pi_{a_{1:i}}$. Our key observation is that
392 from [Equation \(19\)](#), and the convexity of $X \mapsto X^K$, we get

394
$$\frac{q^{\Gamma}(P_1)}{q(P_1)} \leq \frac{q^{\Gamma}(P_2)}{q(P_2)} \leq \dots \leq \frac{q^{\Gamma}(P_M)}{q(P_M)}. \quad (20)$$

396 That is, q^{Γ}/q increases along the global ranking. Hence, to make q^{Γ}/p heuristically close to one,
397 we would also like p/q to increase along the global ranking. Again, enforcing this strict global
398 requirement is difficult, as it requires knowledge of the full joint p_L and q_L . Therefore, we further
399 relax this to a tree-based ranking, by making each local ordering follow p/q in increasing order:

400
$$\pi_{a_{1:i}}(x) = \frac{p(x|a_{1:i})}{q(x|a_{1:i})}. \quad (21)$$

402 We denote the path selection rule induced by these $\pi_{a_{1:i}}$ as Γ_g . Now, we can efficiently compute each
403 $q^{\Gamma_g}(a_{1:i})$ using [Equation \(19\)](#). By definition of the lexicographic ordering, the second sum consists
404 of all paths $b_{1:L}$ where for some $0 \leq j \leq i-1$, we have $b_{1:j} = a_{1:j}$ and $\pi_{a_{1:j}}(b_{j+1}) < \pi_{a_{1:j}}(a_{j+1})$.
405 The first sum contains all these paths, as well as all paths $b_{1:L}$ with $b_{1:i} = a_{1:i}$. This leads to the
406 following closed form expression for q^{Γ_g} :

407
$$q^{\Gamma_g}(a_{1:i}) = \left(q(a_{1:i}) + \sum_{j=0}^{i-1} q(a_{1:j}) \sum_{\pi_{a_{1:j}}(t) < \pi_{a_{1:j}}(a_{j+1})} q(t) \right)^K \quad (22)$$

411
$$- \left(\sum_{j=0}^{i-1} q(a_{1:j}) \sum_{\pi_{a_{1:j}}(t) < \pi_{a_{1:j}}(a_{j+1})} q(t) \right)^K. \quad (23)$$

414 In fact, given a fixed $a_{1:i-1}$, we can efficiently compute all $q^{\Gamma_g}(a_{1:i})$ for $a_i \in \mathcal{V}$. Using the above
415 expression, these $|\mathcal{V}|$ quantities have the exact same terms except at $q(a_{1:i})$ and the $j = i-1$ term
416 in the summation. We can compute the former over all $a_i \in \mathcal{V}$ in $O(|\mathcal{V}|)$ time by multiplying the
417 conditional distribution $q(\cdot|a_{1:i-1})$ by $q(a_{1:i-1})$. The latter is a sum of $q(t)$ over all $t \in \mathcal{V}$ with
418 $\pi_{a_{1:i-1}}(t) < \pi_{a_{1:i-1}}(a_i)$, so we can use a cumulative sum over the ordering on \mathcal{V} induced by $\pi_{a_{1:i-1}}$.
419 From these $|\mathcal{V}|$ values, we can compute conditional probabilities $q^{\Gamma_g}(\cdot|a_{1:i})$ using [Equation \(14\)](#).420 Finally, now that we have an explicit path selection rule Γ_g and corresponding skewed draft q^{Γ_g}
421 values, both of which are efficient to compute and only depend on on-path probabilities, we can
422 follow [Lemma 4.4](#) and use block verification with q^{Γ_g} for the valid single-path verification algorithm.
423 This results in **greedy multi-path block verification (GBV)**, which we explicitly enumerate in
424 [Algorithm 1](#). In the case $K = 1$, this is equivalent to single-path block verification with draft
425 distribution q and target distribution p , because Γ_g only has one path to select, and thus $q^{\Gamma_g} = q$.

426

427 **6 EXPERIMENTS**

428

429 We now test GBV for $K = 1, 2, 3, 4$ paths². As mentioned in [Section 5](#), the case $K = 1$ is equivalent
430 to single-path block verification from [Sun et al. \(2024b\)](#), so it is our baseline approach. For all

431

2We release our code [here](#) with implementations of other multi-path methods.

432 experiments, we select a block length of $L = 8$. We evaluate single-batch sampling from two
 433 target-draft model pairs: OPT 6.7B/350M and OPT 6.7B/125M (Zhang et al., 2022). All temperatures
 434 are set to 1.0. For each pair, we evaluate our algorithm on three datasets: GSM8K, HumanEval, and
 435 MATH500 (Chen et al., 2021; Hendrycks et al., 2021; Cobbe et al., 2021). We take 500 random
 436 problems from the test split of each dataset (HumanEval only has 164). We run our experiments on a
 437 Paperspace machine with an A100-80GB and an Intel Xeon Gold 6342 CPU (12 cores, 2.80 GHz).
 438 We measure *block efficiency* as the number of generated tokens per call to the target model (higher is
 439 better), and *walltime* as average milliseconds per token generated (lower is better).

Model pair	Dataset	Block efficiency (tokens/ M_p -call)			
		$K = 1$	$K = 2$	$K = 3$	$K = 4$
OPT 6.7B/350M	GSM8K	3.294	3.827	3.867	3.884
	HumanEval	3.157	3.538	3.809	3.863
	MATH500	3.227	3.721	3.856	3.896
OPT 6.7B/125M	GSM8K	2.937	3.407	3.526	3.562
	HumanEval	2.653	3.023	3.233	3.503
	MATH500	2.814	3.274	3.392	3.493

450 Table 1: We compute the block efficiency (tokens per target model call) for decoding with GBV, for
 451 target-draft model pairs OPT 6.7B/350M and OPT 6.7B/125M on up to 500 prompts from each of the
 452 datasets GSM8K, HumanEval, and MATH500. Larger numbers are better. In all settings, efficiency
 453 improves as K increases from 1 (block verification baseline) to 4.

455 Our block efficiency results are shown in Table 1. Across all model pairs and datasets, increasing the
 456 number of draft paths K monotonically improves block efficiency. From $K = 1$ to $K = 4$, block
 457 efficiency gains range from 17.91% to 32.04%, with an average 23.08% gain across all six model pair
 458 and dataset combinations. However, we also observe diminishing returns for higher K : on average,
 459 there is a 14.98% increase from $K = 1$ to $K = 2$, a 4.40% increase from $K = 2$ to $K = 3$, and a
 460 2.54% increase from $K = 3$ to $K = 4$. Thus, while GBV significantly improves block efficiency, the
 461 most impactful gains come in the $K = 2$ setting.

Model pair	Dataset	Walltime (ms/token)			
		$K = 1$	$K = 2$	$K = 3$	$K = 4$
OPT 6.7B/350M	GSM8K	44.482	39.592	38.862	48.473
	HumanEval	46.955	43.373	41.017	50.656
	MATH500	45.851	40.473	38.889	48.440
OPT 6.7B/125M	GSM8K	34.479	30.862	30.109	39.831
	HumanEval	38.824	35.648	33.627	41.297
	MATH500	36.022	32.166	31.163	40.329

472 Table 2: We compute the walltimes (ms/token) for decoding with GBV, for target-draft model pairs
 473 OPT 6.7B/350M and OPT 6.7B/125M on up to 500 prompts from each of the datasets GSM8K,
 474 HumanEval, and MATH500. Smaller numbers are better. In all settings, walltimes improve as K
 475 increases from 1 (block verification baseline) to 3, but drop off at $K = 4$.

477 While block efficiency improvements are significant, they do not perfectly align with our walltime
 478 results in Table 2. Here, gains are no longer monotonic from $K = 1$ to $K = 4$. There are significant
 479 improvements from $K = 1$ to $K = 3$, with an average reduction in walltime by 13.34% across
 480 all settings. For OPT 6.7B/350M and MATH500, this even reaches a 15.19% reduction. However,
 481 $K = 4$ always performs worse than the baseline $K = 1$, with an average increase of 9.39% in
 482 walltime. This is due to the increased computation cost of performing the batched target forward pass
 483 over $K = 4$ paths (see Section 2.2), which negates block efficiency gains.

484 We also analyze decoding efficiency across datasets. The relative gains in block efficiency from
 485 $K = 1$ to $K = 4$ are highest for HumanEval (27.20% average across model pairs), followed by
 MATH500 (22.43%) and GSM8K (19.60%). For walltimes, the best setting $K = 3$ reduces walltime

relative to the baseline $K = 1$ by 12.65% for GSM8K, 14.34% for MATH500, and 13.02% for HumanEval (averaged across model pairs). While HumanEval benefits most from extra paths in block efficiency, these do not directly translate to walltime gains.

Furthermore, we examine the impact of target and draft model sizes. Averaged across datasets, the larger draft (350M) achieves 13.72% higher block efficiency than the smaller draft (125M) for $K = 3$. However, the smaller draft achieves a 20.14% lower average walltime than the larger draft in the same setting. While a larger draft model can improve acceptance rates and block efficiency, this comes at the expense of increased latency in draft sequence generation.

For practical usage, we recommend the setting $L = 8, K = 3$. This achieves nearly all of the block efficiency gains of $K = 4$, without incurring a significant spike in latency due to batched target calls. In all settings, $K = 3$ presents the fastest walltime. In multi-batch settings, $K = 2$ is also a viable alternative, as it has around a 10% average walltime reduction compared to block verification. Our results demonstrate that GBV significantly reduces end-to-end latency in autoregressive decoding.

6.1 OTHER MODEL FAMILIES AND DATASETS

While $K = 4$ has the highest block efficiency in our OPT experiments, this is not universally true. In [Appendix I](#), we widen our experimental coverage to include model families with much larger target models (Qwen-3 32B/0.6B and Llama-3 70B/8B) and non-academic datasets (MGSM and ToolBench). In all these settings, we find that block verification outperforms GBV with $K > 1$ in both block efficiency and walltime. We also provide a systems breakdown of relative draft and target pass costs and the KV-cache footprint to explain cost breakdowns for larger models. These results show how to best deploy GBV. For some families (OPT) $K = 3$ is optimal, whereas for others (Qwen-3 and Llama-3) the best choice is to revert to block verification, i.e. GBV with $K = 1$.

6.2 TEMPERATURE AND DRAFT LENGTH ABLATIONS

In [Appendix J](#), we expand our Qwen-3 32B/0.6B experiments for GSM8K in [Appendix I](#) from temperature 1.0 target sampling to 0.2, 0.4, 0.6, and 0.8. We find that block verification performs worse and GBV with $K > 1$ performs better at lower temperatures. At temperature 0.2, GBV with $K = 3$ is over 0.6 tokens/s faster than block verification. Thus, even for model families where block verification is better at temperature one, GBV with $K > 1$ can be better at low temperatures.

In [Appendix K](#), we perform ablations on the draft length L . We find that increasing our $L = 8$ setting to $L = 16, 24$ yields modest gains in block efficiency, but throughput degrades rapidly. Following a systems breakdown similar to [Appendix I](#), we explain this occurs because longer block lengths push more relative work into the drafting bottleneck. We find the ideal choice of L is somewhere around 8.

6.3 COMPARISON TO MULTI-PATH METHODS

Finally, we compare GBV to multi-path verification algorithms SpecTr (Sun et al., 2023), SpecInfer (Miao et al., 2024), and traversal verification (Weng et al., 2025). The rationale for using these methods along with detailed results and discussion are in [Appendix L](#). Traversal verification beats GBV at temperature 1.0 by up to 1.0 tokens/s, GBV beats traversal verification at temperature 0.2 by up to 1.6 tokens/s, and other methods lag behind. In its best setting of $K = 2$, temperature 0.2, GBV achieves nearly 11.5 tokens/s, whereas no other method achieves over 10 tokens/s in *any* setting.

7 CONCLUSION

We developed a single-path information-agnostic linear program (LP) which encodes feasible speedups for valid verification algorithms in speculative sampling, and showed that block verification remains optimal even with access to off-path probabilities. We further extended our LP to the setting where multiple draft paths are generated. While this multi-path LP is not feasible to solve exactly, by approximating it with a class of greedy verification algorithms, we developed a generalization of block verification, called greedy multi-path block verification. This significantly improves decoding efficiency relative to block verification. Future work could further explore the class of greedy verification schemes, or explore alternative approximation to the multi-path LP.

540 REPRODUCIBILITY STATEMENT
541542 The majority of our work is theoretical, with proofs in Appendices A to G. While we do not release
543 code, we provide enough detail in [Section 5](#) to reproduce our greedy multi-path block verification
544 algorithm. We also provide details around our machine setup, datasets, and model pairs in [Section 6](#).
545546 REFERENCES
547

548 Amey Agrawal, Nitin Kedia, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav S. Gula-
549 vani, Alexey Tumanov, and Ramachandran Ramjee. Taming throughput-latency tradeoff in llm
550 inference with sarathi-serve. In *18th USENIX Symposium on Operating Systems Design and*
551 *Implementation (OSDI '24)*, 2024. URL <https://www.usenix.org/system/files/osdi24-agrawal.pdf>.

553 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
554 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
555

556 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
557 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models
558 are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp.
559 1877–1901, 2020.

560 Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D Lee, Deming Chen, and Tri Dao.
561 Medusa: Simple llm inference acceleration framework with multiple decoding heads. *arXiv*
562 *preprint arXiv:2401.10774*, 2024.
563

564 Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John
565 Jumper. Accelerating large language model decoding with speculative sampling. *arXiv preprint*
566 *arXiv:2302.01318*, 2023.

567 Mark Chen, Jerry Tworek, Heewoo Jun, et al. Evaluating large language models trained on code.
568 *arXiv preprint arXiv:2107.03374*, 2021.
569

570 Ziyi Chen, Xiaocong Yang, Jiacheng Lin, Chenkai Sun, Kevin Chang, and Jie Huang. Cascade
571 speculative drafting for even faster llm inference. *Advances in Neural Information Processing*
572 *Systems*, 37:86226–86242, 2024.

573 Karl Cobbe, Vineet Kosaraju, Mohammad Ostendorf, et al. Training verifiers to solve math word
574 problems. *arXiv preprint arXiv:2110.14168*, 2021.
575

576 Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast
577 and memory-efficient exact attention with io-awareness. In *Advances in Neural Information*
578 *Processing Systems (NeurIPS 2022)*, 2022. URL <https://openreview.net/pdf?id=H4DqfPSibmx>.

579

580 Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich, Basil Hosmer, Bram Wasti, Liangzhen Lai,
581 Anas Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed Roman, et al. Layerskip: Enabling early
582 exit inference and self-speculative decoding. *arXiv preprint arXiv:2404.16710*, 2024.
583

584 Yichao Fu, Peter Bailis, Ion Stoica, and Hao Zhang. Break the sequential dependency of llm inference
585 using lookahead decoding. *arXiv preprint arXiv:2402.02057*, 2024.
586

587 David Gale. A theorem on flows in networks. In *Classic Papers in Combinatorics*, pp. 259–268.
588 Springer, 1957.

589 Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve.
590 Better & faster large language models via multi-token prediction. *arXiv preprint arXiv:2404.19737*,
591 2024.

592

593 Zhenyu He, Zexuan Zhong, Tianle Cai, Jason D Lee, and Di He. Rest: Retrieval-based speculative
decoding. *arXiv preprint arXiv:2311.08252*, 2023.

594 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
 595 Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint*
 596 *arXiv:2103.03874*, 2021.

597

598 Zhengmian Hu and Heng Huang. Accelerated speculative sampling based on tree monte carlo. In
 599 *Forty-first International Conference on Machine Learning*, 2024.

600 Zhengmian Hu, Tong Zheng, Vignesh Viswanathan, Ziyi Chen, Ryan A Rossi, Yihan Wu, Dinesh
 601 Manocha, and Heng Huang. Towards optimal multi-draft speculative decoding. *arXiv preprint*
 602 *arXiv:2502.18779*, 2025.

603

604 Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank
 605 Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, Stephan Krusche, Gitta
 606 Kutyniok, Tilman Michaeli, Claudia Nerdel, Juergen Pfeffer, Oleksandra Poquet, Michael Sailer,
 607 Albrecht Schmidt, Tina Seidel, and Gjergji Kasneci. Chatgpt for good? on opportunities and
 608 challenges of large language models for education. *Learning and Individual Differences*, 103:
 609 102274, 01 2023. doi: 10.1016/j.lindif.2023.102274.

610

611 Ashish J Khisti, MohammadReza Ebrahimi, Hassan Dbouk, Arash Behboodi, Roland Memisevic,
 612 and Christos Louizos. Multi-draft speculative sampling: Canonical decomposition and theoretical
 613 limits. In *The Thirteenth International Conference on Learning Representations*, 2025.

614

615 Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative
 616 decoding. In *International Conference on Machine Learning*, pp. 19274–19286. PMLR, 2023.

617

618 Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. Eagle: Speculative sampling requires
 619 rethinking feature uncertainty. *arXiv preprint arXiv:2401.15077*, 2024.

620

621 Xiaoxuan Liu, Lanxiang Hu, Peter Bailis, Alvin Cheung, Zhijie Deng, Ion Stoica, and Hao Zhang.
 622 Online speculative decoding. *arXiv preprint arXiv:2310.07177*, 2023.

623

624 Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae
 625 Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, et al. Specinfer: Accelerating large
 626 language model serving with tree-based speculative inference and verification. In *Proceedings of
 627 the 29th ACM International Conference on Architectural Support for Programming Languages
 628 and Operating Systems, Volume 3*, pp. 932–949, 2024.

629

630 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

631

632 Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language un-
 633 derstanding by generative pre-training. [https://cdn.openai.com/research-covers/
 634 language-unsupervised/language_understanding_paper.pdf](https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf), 2018.

635

636 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever.
 637 Language models are unsupervised multitask learners. [https://cdn.openai.
 638 com/better-language-models/language_models_are_unsupervised_](https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf)
 639 [multitask_learners.pdf](https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf), 2019.

640

641 Mohammad Samragh, Arnav Kundu, David Harrison, Kumari Nishu, Devang Naik, Minsik Cho, and
 642 Mehrdad Farajtabar. Your llm knows the future: Uncovering its multi-token prediction potential.
 643 *arXiv preprint arXiv:2507.11851*, 2025.

644

645 Alexander Schrijver et al. *Combinatorial optimization: polyhedra and efficiency*, volume 24. Springer,
 646 2003.

647

648 Benjamin Spector and Chris Re. Accelerating llm inference with staged speculative decoding. *arXiv
 649 preprint arXiv:2308.04623*, 2023.

650

651 Hanshi Sun, Zhuoming Chen, Xinyu Yang, Yuandong Tian, and Beidi Chen. Triforce: Lossless
 652 acceleration of long sequence generation with hierarchical speculative decoding. *arXiv preprint
 653 arXiv:2404.11912*, 2024a.

648 Ziteng Sun, Ananda Theertha Suresh, Jae Hun Ro, Ahmad Beirami, Himanshu Jain, and Felix
 649 Yu. Spectr: Fast speculative decoding via optimal transport. *Advances in Neural Information
 650 Processing Systems*, 36:30222–30242, 2023.

651 Ziteng Sun, Uri Mendlovic, Yaniv Leviathan, Asaf Aharoni, Jae Hun Ro, Ahmad Beirami, and
 652 Ananda Theertha Suresh. Block verification accelerates speculative decoding. *arXiv preprint
 653 arXiv:2403.10444*, 2024b.

654 Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang
 655 Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature Medicine*, 29:
 656 1930–1940, 2023. doi: 10.1038/s41591-023-02459-w. URL <https://www.nature.com/articles/s41591-023-02459-w>. Review Article, Published: 17 July 2023.

657 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 658 Lacroix, Naman Goyal, Eric Hambro, Haoran Azhar, Alice Rodriguez, et al. Llama 2: Open
 659 foundation and fine-tuned chat models. *arXiv preprint arXiv:2310.11387*, 2023a.

660 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 661 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Haoran Azhar, et al. Llama: Open and
 662 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023b.

663 Yepeng Weng, Qiao Hu, Xujie Chen, Li Liu, Dianwen Mei, Huishi Qiu, Jiang Tian, and Zhongchao
 664 Shi. Traversal verification for speculative tree decoding. *arXiv preprint arXiv:2505.12398*, 2025.

665 Jun Zhang, Jue Wang, Huan Li, Lidan Shou, Ke Chen, Gang Chen, and Sharad Mehrotra. Draft &
 666 verify: Lossless large language model acceleration via self-speculative decoding. *arXiv preprint
 667 arXiv:2309.08168*, 2023.

668 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuhui Chen, Christopher
 669 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language
 670 models. *arXiv preprint arXiv:2205.01068*, 2022.

671 Yongchao Zhou, Kaifeng Lyu, Ankit Singh Rawat, Aditya Krishna Menon, Afshin Rostamizadeh,
 672 Sanjiv Kumar, Jean-François Kagy, and Rishabh Agarwal. Distillspec: Improving speculative
 673 decoding via knowledge distillation. *arXiv preprint arXiv:2310.08461*, 2023.

674 Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen,
 675 and Lei Li. Multilingual machine translation with large language models: Empirical results and
 676 analysis, 2024. URL <https://arxiv.org/abs/2304.04675>.

677
 678
 679
 680
 681
 682
 683
 684
 685
 686
 687
 688
 689
 690
 691
 692
 693
 694
 695
 696
 697
 698
 699
 700
 701

702 A PROOF OF MULTI-PATH INFORMATION-AGNOSTIC LP

704 In this section, we prove the multi-path information-agnostic LP in [Theorem 4.3](#). Our proof begins
 705 with a key lemma that dictates when one can sample from one distribution given data from another,
 706 with the conditional distribution (often called the transport) to induce the former restricted to a
 707 given support. In our context, the restricted support represents the prefix-matching requirement, the
 708 distribution from which we are given data is induced by i.i.d. sampling paths from the draft model,
 709 and the distribution from which we would like to sample is that of the verification algorithm output.

710 **Lemma A.1** (Bipartite Transport Feasibility). *Let $G = (A \cup B, E \subseteq A \times B)$ be a bipartite graph,
 711 with probability distributions $a(\cdot)$ and $b(\cdot)$ over A and B , respectively. Then the following conditions
 712 are equivalent, where $N(\cdot)$ denotes neighborhoods:*

- 714 1. *There exists a joint distribution $\pi(\cdot, \cdot)$ over $A \times B$ with marginal distributions $a(\cdot)$ and $b(\cdot)$,
 715 such that $\pi(x, y) = 0$ for all $(x, y) \notin E$.*
- 716 2. *For any $S \subseteq A$, we have $\sum_{x \in S} a(x) \leq \sum_{y \in N(S)} b(y)$.*

719 *Proof.* Condition 1 is equivalent to the feasibility of the following LP in variables $\pi_{x,y} \geq 0$:

$$721 \sum_{x \in A} \pi_{x,y} = b(y) \quad \forall y \in B, \quad \sum_{y \in B} \pi_{x,y} = a(x) \quad \forall x \in A, \quad \pi_{x,y} = 0 \quad \forall (x, y) \notin E. \quad (24)$$

724 We can incorporate the zero equality condition into the sum equalities to turn this into:

$$726 \sum_{x \in N(y)} \pi_{x,y} = b(y) \quad \forall y \in B, \quad \sum_{y \in N(x)} \pi_{x,y} = a(x) \quad \forall x \in A. \quad (25)$$

728 Now, Gale's feasibility theorem for bipartite supply-demand networks ([Gale, 1957](#)) implies this LP is
 729 feasible in nonnegative variables if and only if
 730

$$731 \sum_{x \in S} a(x) \leq \sum_{y \in N(S)} b(x) \quad \forall S \subseteq A, \quad \sum_{x \in A} a(x) = \sum_{y \in B} b(x). \quad (26)$$

734 The equality condition is automatically satisfied, since $a(\cdot), b(\cdot)$ are probability distributions, so both
 735 sides become one. Thus, Condition 1 is equivalent to Condition 2. \square

737 Using this, we can prove the following lemma, which dictates what distributions are feasible for a
 738 K -path verification algorithm Φ which satisfies prefix-matching. We do not enforce target-matching
 739 yet, so this algorithm is not necessarily valid. For any K -path verification algorithm Φ , denote $R_\Phi(\cdot)$
 740 as the distribution of the output of Φ , which is supported over $\mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$, i.e. all nonempty paths
 741 of length $\leq L+1$. We need the nonempty requirement as we must output at least one token. As we
 742 prove in [Lemma A.2](#), a collection of subset-sum inequalities over a particular bipartite graph exactly
 743 describes the set of R_Φ which can be realized by some Φ satisfying prefix-matching.

744 **Lemma A.2** (Prefix-Matching). *Construct the bipartite graph G with left vertices $\mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$, right
 745 vertices $(\mathcal{V}^L)^K$, and edges E consisting of all pairs*

$$747 \left(a_{1:i}, \left(a_{1:L}^{(1)}, \dots, a_{1:L}^{(K)} \right) \right) \in (\mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0) \times (\mathcal{V}^L)^K, \quad a_{1:i-1} \in \left\{ a_{1:i-1}^{(1)}, \dots, a_{1:i-1}^{(K)} \right\}. \quad (27)$$

749 The values $R_\Phi(a_{1:i})$ for $a_{1:i} \in \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$ are feasible for a K -path verification algorithm Φ that
 750 satisfies prefix-matching if and only if they are nonnegative, sum to one, and

$$752 \sum_{a_{1:i} \in S} R_\Phi(a_{1:i}) \leq \sum_{\left(a_{1:L}^{(1)}, \dots, a_{1:L}^{(K)} \right) \in N(S)} \prod_{k=1}^K q \left(a_{1:L}^{(k)} \right) \quad \forall S \subseteq \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0, \quad (28)$$

755 where $N(\cdot)$ denotes a neighborhood in G .

756 *Proof.* The sum to one equality and nonnegative requirements are equivalent to requiring that R_Φ is
 757 a valid probability distribution. Now, given that R_Φ is a valid probability distribution, we observe
 758 that the prefix-matching requirement is equivalent to the following: given K i.i.d. sampled paths
 759

$$760 \quad X_{1:L}^{(1)}, \dots, X_{1:L}^{(K)} \sim q(\cdot), \quad (29)$$

761 there is a conditional distribution (randomized rule) $\pi(\cdot | X_{1:L}^{(1)}, \dots, X_{1:L}^{(K)})$ over $\mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$ which
 762 is supported only over prefixes of these paths plus an additional token, such that sampling from π
 763 results in the distribution R_Φ . Given the construction of the graph G , we can alternatively express
 764 the restricted support requirement on $\pi(\cdot | \cdot)$ as $\pi(v_l | v_r) = 0$ for all left vertices v_l and right vertices
 765 v_r with $(v_l, v_r) \notin E$. Thus, by turning this conditional distribution into a joint distribution, as R_Φ
 766 is a distribution over left vertices and i.i.d. sampling from $q(\cdot)$ is a distribution over right vertices,
 767 [Lemma A.1](#) directly applies to show that prefix-matching is equivalent to
 768

$$769 \quad \sum_{a_{1:i} \in S} R_\Phi(a_{1:i}) \leq \sum_{\substack{(a_{1:L}^{(1)}, \dots, a_{1:L}^{(K)}) \in N(S)}} \prod_{k=1}^K q(a_{1:L}^{(k)}) \quad \forall S \subseteq \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0. \quad (30)$$

772 The product term is the probability of sampling K paths i.i.d from q . This completes the proof. \square
 773

774 Now, we move onto the target-matching requirement. Again, we can establish a necessary and
 775 sufficient condition for a K -path verification algorithm Φ to satisfy target-matching terms of the R_Φ
 776 values. Now, the condition is path-wise summation requirement.

777 **Lemma A.3** (Target-Matching). *The values $R_\Phi(a_{1:i})$ for $a_{1:i} \in \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$ are feasible for a
 778 K -path verification algorithm Φ that satisfies prefix-matching if and only if they are nonnegative,
 779 sum to one, and*

$$780 \quad \sum_{i=0}^L \frac{R_\Phi(a_{1:i+1})}{p(a_{1:i+1})} = 1 \quad \forall a_{1:L+1} \in \mathcal{V}^{L+1}. \quad (31)$$

784 *Proof.* Target-matching asserts that sampling from $p_{L+1}(\cdot)$ is equivalent to sampling the output
 785 ($X_{1:\tau}, Y$) $\sim R_\Phi$ of the verification algorithm and then sampling an additional $L - \tau$ tokens from
 786 $Z_{1:L-\tau} \sim p_{L-\tau}(\cdot | X_{1:\tau}, Y)$. For a given path $a_{1:L+1} \in \mathcal{V}^{L+1}$, the probability of sampling this path
 787 from the former method is $p_{L+1}(a_{1:L+1})$. The probability of sampling it from the latter method is
 788

$$789 \quad \mathbb{P}(X_{1:\tau} = a_{1:\tau}, Y = a_{\tau+1}, Z_{1:L-\tau} = a_{\tau+2:L+1}) \quad (32)$$

$$790 \quad = \sum_{i=0}^L \mathbb{P}(\tau = i, X_{1:\tau} = a_{1:\tau}, Y = a_{\tau+1}, Z_{1:L-\tau} = a_{\tau+2:L+1}) \quad (33)$$

$$791 \quad = \sum_{i=0}^L \mathbb{P}(\tau = i, X_{1:\tau} = a_{1:i}, Y = a_{i+1}, Z_{1:L-\tau} = a_{i+2:L+1}) \quad (34)$$

$$792 \quad = \sum_{i=0}^L \mathbb{P}((X_{1:\tau}, Y) = a_{1:i+1}) p(a_{i+2:L+1} | a_{1:i+1}) \quad (35)$$

$$793 \quad = \sum_{i=0}^L R_\Phi(a_{1:i+1}) \cdot \frac{p(a_{1:L+1})}{p(a_{1:i+1})}. \quad (36)$$

802 Setting these two equal for all paths $a_{1:L+1}$ and dividing out the $p(a_{1:L+1}) = p_{L+1}(a_{1:L+1})$ term
 803 gives the desired condition. \square

804 By combining the prefix-matching and target-matching lemmas, we can now write down the multi-
 805 path LP, in a form which is slightly more complicated than in [Theorem 4.3](#), and is in terms of variables
 806 R_Φ rather than the node budgets D_Φ . We will use the notion of a parent set and minimal set.

807 **Definition A.4.** For each $S \subseteq \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$, the **parent set** $\text{par}(S) \subseteq \mathcal{V}^{\leq L}$ of S consists of all
 808 $a_{1:i} \in \mathcal{V}^{\leq L}$ where $a_{1:i+1} \in S$ for some $a_{i+1} \in \mathcal{V}$. For each $S \subseteq \mathcal{V}^{\leq L}$, the **minimal set** $\text{min}(S) \subseteq S$
 809 of S consists of all $a_{1:i} \in S$ where no proper prefix of $a_{1:i}$ appears in S , i.e. $a_{1:0}, \dots, a_{1:i-1} \notin S$.

Note that an antichain (Definition 4.2) is precisely a subset $S \subseteq \mathcal{V}^{\leq L}$ where $\min(S) = S$, and any minimal set $\min(S)$ is also an antichain because $\min(\min(S)) = \min(S)$. We will use parent and minimal sets to remove redundant inequalities in the Lemma A.2.

Lemma A.5 (Unsimplified Multi-Path LP). *The values $R_\Phi(a_{1:i})$ for $a_{1:i} \in \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$ are feasible for a valid K -path verification algorithm Φ if and only if they are nonnegative, sum to one, and are feasible in*

$$\max \sum_{i=1}^{L+1} \sum_{a_{1:i} \in \mathcal{V}^i} i R_\Phi(a_{1:i}), \quad (37)$$

$$\text{s.t. } \sum_{i=0}^L \frac{R_\Phi(a_{1:i+1})}{p(a_{1:i+1})} = 1 \quad \forall a_{1:L+1} \in \mathcal{V}^{L+1}, \quad (38)$$

$$\sum_{a_{1:i} \in T} \sum_{j=i+1}^{L+1} \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} R_\Phi(a_{1:j}) \leq 1 - \left(1 - \sum_{a_{1:i} \in T} q(a_{1:i})\right)^K \quad \forall T \in \mathcal{A}(\mathcal{V}^{\leq L}). \quad (39)$$

Furthermore, the objective value is precisely the block efficiency $\mathbb{E}[\tau + 1]$ for Φ .

Proof. First, fix a subset $S \subseteq \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$. The neighborhood $N(S) \subseteq (\mathcal{V}^L)^K$ induced by the graph G in Lemma A.2 is the set of K -tuples of length- L paths where at least one path has a prefix in $\text{par}(S)$. So, the complement of $N(S)$ consists of all K -tuples of length- L paths where no path has a prefix in $\text{par}(S)$. Now, note that a path has a prefix in $\text{par}(S)$ if and only if it has a prefix in $\min(\text{par}(S))$. Because $\min(\text{par}(S))$ is an antichain (i.e. no two distinct elements in it are a prefix of the same path), disjointness shows that the probability that a length L -path sampled from $q(\cdot)$ has a prefix in $\min(\text{par}(S))$ is $\sum_{a_{1:i} \in \min(\text{par}(S))} q(a_{1:i})$. Using these observations, we simplify:

$$\sum_{(a_{1:L}^{(1)}, \dots, a_{1:L}^{(K)}) \in N(S)} \prod_{k=1}^K q(a_{1:L}^{(k)}) = 1 - \sum_{(a_{1:L}^{(1)}, \dots, a_{1:L}^{(K)}) \notin N(S)} \prod_{k=1}^K q(a_{1:L}^{(k)}) \quad (40)$$

$$= 1 - \left(1 - \sum_{a_{1:i} \in \min(\text{par}(S))} q(a_{1:i})\right)^K. \quad (41)$$

Thus, the prefix-matching inequalities from Lemma A.2 become:

$$\sum_{a_{1:i} \in S} R_\Phi(a_{1:i}) \leq 1 - \left(1 - \sum_{a_{1:i} \in \min(\text{par}(S))} q(a_{1:i})\right)^K \quad \forall S \subseteq \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0. \quad (42)$$

We now make the key observation that many of these inequalities are redundant. For any $a_{1:i}$ with a proper (i.e. not full) prefix in $\min(\text{par}(S))$, we can add $a_{1:i}$ to S without changing $\min(\text{par}(S))$, keeping the right hand side constant. This also increases the left hand side (since R_Φ is nonnegative). Thus, the strictest of these inequalities are precisely those where S includes all paths with a proper prefix in $\min(\text{par}(S))$. For a fixed $T = \min(\text{par}(S))$, this makes S the set of all $a_{1:j}$ where $a_{1:i} \in T$ for some $i < j$. Since all antichains are achievable as some $\min(\text{par}(S))$, this turns the above collection of inequalities into

$$\sum_{a_{1:i} \in T} \sum_{j=i+1}^{L+1} \sum_{a_{i+1:j} \in \mathcal{V}^{\leq j-i}} R_\Phi(a_{1:j}) \leq 1 - \left(1 - \sum_{a_{1:i} \in T} q(a_{1:i})\right)^K \quad \forall T \in \mathcal{A}(\mathcal{V}^{\leq L}). \quad (43)$$

This covers the prefix-matching. The target-matching requirement remains the same as in Lemma A.3. Thus, all that remains is to show the objective values is the block efficiency. This holds as $R_\Phi(a_{1:i})$ is the probability that Φ outputs $a_{1:i}$, and the total number of tokens generated with this output is i , so summing over all possible outputs $a_{1:i} \in \mathcal{V}^{\leq L+1} \setminus \mathcal{V}^0$ gives the desired result. \square

Finally, we can prove the multi-path LP in Theorem 4.3. The only remaining step to derive this is to turn the variables R_Φ into variables D_Φ representing node budgets.

864 *Proof.* By the definition of D_Φ in [Definition 3.2](#), we see that
865

$$866 \quad 867 \quad 868 \quad D_\Phi(a_{1:i}) = 1 - \sum_{j=1}^i \frac{R_\Phi(a_{1:i})}{p(a_{1:i})} \quad \forall a_{1:i} \in \mathcal{V}^{\leq L+1}. \quad (44)$$

869 The target-matching condition simply implies that each $D_\Phi(a_{1:L+1}) = 0$. Now, observe that
870

$$871 \quad R_\Phi(a_{1:i}) = D_\Phi(a_{1:i-1})p(a_{1:i}) - D_\Phi(a_{1:i})p(a_{1:i}) \quad \forall a_{1:i} \in \mathcal{V}^{\leq L}. \quad (45)$$

872 for all. Thus, incorporating the nonnegativity constraint on R_Φ , we get
873

$$874 \quad 1 = D_\Phi(a_{1:0}) \geq D_\Phi(a_{1:1}) \geq \dots \geq D_\Phi(a_{1:L}) \geq D_\Phi(a_{1:L+1}) = 0 \quad \forall a_{1:L+1} \in \mathcal{V}^{L+1}. \quad (46)$$

875 Next, observe that for any $a_{1:i} \in \mathcal{V}^{\leq L}$, we have the following telescoping identity:
876

$$877 \quad \sum_{j=i+1}^{L+1} \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} R_\Phi(a_{1:j}) \quad (47)$$

$$878 \quad = \sum_{j=i+1}^{L+1} \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} D_\Phi(a_{1:j-1})p(a_{1:j}) - \sum_{j=i+1}^{L+1} \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} D_\Phi(a_{1:j})p(a_{1:j}) \quad (48)$$

$$879 \quad = \sum_{j=i}^L \sum_{a_{i+1:j+1} \in \mathcal{V}^{j+1-i}} D_\Phi(a_{1:j})p(a_{1:j+1}) - \sum_{j=i+1}^{L+1} \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} D_\Phi(a_{1:j})p(a_{1:j}) \quad (49)$$

$$880 \quad = D_\Phi(a_{1:i})p(a_{1:i}) + \sum_{j=i+1}^L \sum_{a_{i+1:j+1} \in \mathcal{V}^{j+1-i}} D_\Phi(a_{1:j})p(a_{1:j+1}) \quad (50)$$

$$881 \quad - \sum_{j=i+1}^L \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} D_\Phi(a_{1:j})p(a_{1:j}) - 0 \quad (51)$$

$$882 \quad = D_\Phi(a_{1:i})p(a_{1:i}) + \sum_{j=i+1}^L D_\Phi(a_{1:j}) \left(\sum_{a_{i+1:j+1} \in \mathcal{V}^{j+1-i}} p(a_{1:j+1}) - \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} p(a_{1:j}) \right) \quad (52)$$

$$883 \quad = D_\Phi(a_{1:i})p(a_{1:i}) + \sum_{j=i+1}^L D_\Phi(a_{1:j}) (p(a_{1:i}) - p(a_{1:j})) = D_\Phi(a_{1:i})p(a_{1:i}). \quad (53)$$

884 This means the condition that all R_Φ sum to one is automatically satisfied by applying this identity
885 for $i = 0$. This also reduces the antichain condition in the unsimplified multi-path LP to
886

$$887 \quad \sum_{a_{1:i} \in T} D_\Phi(a_{1:i})p(a_{1:i}) \leq 1 - \left(1 - \sum_{a_{1:i} \in T} q(a_{1:i}) \right)^K \quad \forall T \in \mathcal{A}(\mathcal{V}^{\leq L}). \quad (54)$$

888 This covers all conditions, so all that remains is to show that the objective reduces to the desired
889 expression. This holds by interchanging the order of summation, and using the telescoping identity:
890

$$891 \quad \sum_{i=1}^{L+1} \sum_{a_{1:i} \in \mathcal{V}^i} i R_\Phi(a_{1:i}) = \sum_{i=1}^{L+1} \sum_{j=0}^{i-1} \sum_{a_{1:i} \in \mathcal{V}^i} R_\Phi(a_{1:i}) \quad (55)$$

$$892 \quad = \sum_{j=0}^L \sum_{i=j+1}^{L+1} \sum_{a_{1:i} \in \mathcal{V}^i} R_\Phi(a_{1:i}) \quad (56)$$

$$893 \quad = \sum_{j=0}^L \sum_{i=j+1}^{L+1} \sum_{a_{1:i} \in \mathcal{V}^j} \sum_{a_{j+1:i} \in \mathcal{V}^{j-i}} R_\Phi(a_{1:i}) \quad (57)$$

$$894 \quad = \sum_{j=0}^L \sum_{a_{1:j} \in \mathcal{V}^j} \sum_{i=j+1}^{L+1} \sum_{a_{j+1:i} \in \mathcal{V}^{j-i}} R_\Phi(a_{1:i}) = \sum_{a_{1:j} \in \mathcal{V}^{\leq L}} D_\Phi(a_{1:j})p(a_{1:j}). \quad (58)$$

910 This completes the reduction to the desired form for the multi-path LP. \square

918 B PROOF OF SINGLE-PATH INFORMATION-AGNOSTIC LP

920 Using the multi-path in [Theorem 4.3](#), it is easy to prove the single-path LP in [Theorem 3.3](#) by reducing
 921 it to the case $K = 1$. This works because a valid single-path verification algorithm is precisely the
 922 same as a valid K -path verification algorithm in the case $K = 1$.

924 *Proof.* It suffices to show that when $K = 1$, the multi-path information-agnostic LP from [Theorem 4.3](#)
 925 reduces to the single-path LP. Indeed, the antichain condition becomes

$$926 \quad \sum_{a_{1:i} \in T} D_\Phi(a_{1:i}) p(a_{1:i}) \leq \sum_{a_{1:i} \in T} q(a_{1:i}) \quad \forall T \in \mathcal{A}(\mathcal{V}^{\leq L}). \quad (59)$$

929 This forces the pointwise conditions $D_\Phi(a_{1:i}) p(a_{1:i}) \leq q(a_{1:i})$, since by taking the singleton
 930 antichains $T = \{a_{1:i}\}$. Furthermore, these pointwise conditions imply the antichain conditions
 931 through summation. All other conditions and the objective function remain the same. Hence, we
 932 obtain the desired single-path LP. \square

933 C PROOF OF BLOCK VERIFICATION OPTIMALITY

935 Using the single-path LP, we can prove that block verification is optimal in the class of valid single-
 936 path verification algorithms ([Theorem 3.4](#)). We will use results from [Sun et al. \(2024b\)](#). The idea is to
 937 show that $D_{\Phi_{BV}}$ for the single-path block verification algorithm Φ_{BV} are tight in the single-path LP.
 938

939 *Proof.* From Lemma 4 in Appendix B of [Sun et al. \(2024b\)](#), we have for each path $a_{1:i} \in \mathcal{V}^{\leq L}$ that

$$940 \quad \mathbb{P}(\tau \geq i | X_{1:i} = a_{1:i}) = w(a_{1:i}) \quad (60)$$

942 under Φ_{BV} , where weights w are defined in [Equation \(3\)](#). Now, because Φ_{BV} is a valid verification
 943 algorithm, we can use the telescoping identity from the proof of [Theorem 4.3](#) at the end of [Appendix A](#)
 944 (with $R_{\Phi_{BV}}$ defined similarly as the distribution of the output of Φ_{BV}) to show

$$945 \quad D_{\Phi_{BV}}(a_{1:i}) p(a_{1:i}) = \sum_{j=i+1}^{L+1} \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} R_{\Phi_{BV}}(a_{1:j}) \quad (61)$$

$$948 \quad = \sum_{j=i+1}^{L+1} \sum_{a_{i+1:j} \in \mathcal{V}^{j-i}} \mathbb{P}(\tau = j-1, X_{1:j-1} = a_{1:j-1}, Y = a_j) \quad (62)$$

$$951 \quad = \sum_{j=i+1}^{L+1} \mathbb{P}(\tau = j-1, X_{1:i} = a_{1:i}) = \mathbb{P}(\tau \geq i, X_{1:i} = a_{1:i}) \quad (63)$$

$$954 \quad = \mathbb{P}(\tau \geq i | X_{1:i} = a_{1:i}) \mathbb{P}(X_{1:i} = a_{1:i}) = w(a_{1:i}) q(a_{1:i}). \quad (64)$$

955 Thus, to show block verification is optimal, it suffices to show that for any valid single-path verification
 956 algorithm Φ and path $a_{1:i} \in \mathcal{V}^{\leq L}$, we have $D_\Phi(a_{1:i}) p(a_{1:i}) \leq w(a_{1:i}) q(a_{1:i})$, as then the objective
 957 value for Φ (the block efficiency, i.e. the sum of these $D_\Phi p$ terms) cannot exceed the objective value
 958 for Φ_{BV} (the sum of the $D_{\Phi_{BV}} p$ terms) in the single-path LP. The proof is by induction on i . For the
 959 base case $i = 0$, this is clear as $w(a_{1:0}) q(a_{1:0}) = 1 = D_\Phi(a_{1:0}) p(a_{1:0})$. Now, suppose this statement
 960 holds for i . To prove it for $i + 1$, we use the single-path conditions to obtain

$$961 \quad D_\Phi(a_{1:i+1}) p(a_{1:i+1}) \leq D_\Phi(a_{1:i}) p(a_{1:i+1}) \quad (65)$$

$$962 \quad = D_\Phi(a_{1:i}) p(a_{1:i}) p(a_{i+1} | a_{1:i}) \quad (66)$$

$$963 \quad \leq w(a_{1:i}) q(a_{1:i}) p(a_{i+1} | a_{1:i}) \quad (67)$$

$$965 \quad = q(a_{1:i+1}) \cdot \frac{p(a_{i+1} | a_{1:i})}{q(a_{i+1} | a_{1:i})} \cdot w(a_{1:i}) \quad (68)$$

$$967 \quad \leq q(a_{1:i+1}) \cdot \min \left\{ 1, \frac{p(a_{i+1} | a_{1:i})}{q(a_{i+1} | a_{1:i})} \cdot w(a_{1:i}) \right\} = q(a_{1:i+1}) w(a_{1:i+1}). \quad (69)$$

968 The last inequality holds by using the pointwise $D_\Phi p \leq q$ bound from the single-path LP. This
 969 completes the induction, and thus the proof that block verification is optimal.

971 \square

972 **D PROOF OF PATH SELECTION DECOMPOSITION**
973

974 In this section, we prove [Lemma 4.4](#), which states that any valid K -path verification algorithm can
975 be decomposed into a randomized path selection rule which chooses one of the K paths, and then a
976 valid single-path verification algorithm that operates on an skewed draft distribution.
977

978 *Proof.* Let Φ denote the valid verification algorithm. This takes in K paths P_1, \dots, P_K i.i.d. sampled
979 from $q(\cdot)$, and samples a nonempty path in $\mathcal{V}^{\leq L+1}$ using a conditional distribution $\pi(\cdot | P_1, \dots, P_K)$,
980 which can depend on off-path target and draft distribution values. This is constrained to only output
981 $(X_{1:\tau}, Y)$ where $X_{1:\tau}$ is a prefix of one of P_1, \dots, P_K and $Y \in \mathcal{V}$. Thus, we can conditionally
982 factorize π into a prefix selection rule π_1 on prefixes of P_1, \dots, P_K conditioned over P_1, \dots, P_K ,
983 and an additional token rule π_2 on \mathcal{V} conditioned over $P_1, \dots, P_K, X_{1:\tau}$:
984

985
$$\pi(X_{1:\tau}, Y | P_1, \dots, P_K) = \pi_1(X_{1:\tau} | P_1, \dots, P_K) \pi_2(Y | P_1, \dots, P_K, X_{1:\tau}). \quad (70)$$

986

987 For each prefix $a_{1:i}$ of some path in P_1, \dots, P_K , denote its multiplicity $m(a_{1:i}) \geq 1$ as the number
988 of paths that contain it as a prefix. For example, $m(\emptyset) = K$. Then, we define a path selection rule π_3
989 on $[K]$ conditioned over P_1, \dots, P_K , as:

990
$$\pi_3(k_0 | P_1, \dots, P_K) = \sum_{i=0}^L \pi_1((P_{k_0})_{1:i} | P_1, \dots, P_K) m((P_{k_0})_{1:i})^{-1}. \quad (71)$$

991

992 To see this is a valid probability distribution, observe that for each distinct $a_{1:i}$ appearing as a prefix
993 of one of P_1, \dots, P_K , when we sum the above formula over all $k_0 \in [K]$, we have exactly $m(a_{1:i})$
994 terms in the sum above with $(P_{k_0})_{1:i} = a_{1:i}$, each of which are $\pi_1(a_{1:i} | P_1, \dots, P_K) m(a_{1:i})^{-1}$.
995 Thus, the terms containing $a_{1:i}$ cancel to $\pi_1(a_{1:i} | P_1, \dots, P_K)$, and summing this over all *distinct*
996 $a_{1:i}$ gives a result of 1, because π_1 is a valid distribution. Also, define a prefix selection rule π_4 on
997 prefixes of P_{k_0} conditioned over $P_1, \dots, P_K, k_0 \in [K]$:
998

999
$$\pi_4((P_{k_0})_{1:i} | P_1, \dots, P_K, k_0) = \frac{\pi_1((P_{k_0})_{1:i} | P_1, \dots, P_K) m((P_{k_0})_{1:i})^{-1}}{\sum_{i=0}^L \pi_1((P_k)_{1:i} | P_1, \dots, P_K) m((P_k)_{1:i})^{-1}}, \quad (72)$$

1000

1001 which is also a valid probability distribution. We see that when sampling $k_0 \sim \pi_3(\cdot | P_1, \dots, P_K)$ and
1002 then $X_{1:\tau} \sim \pi_4(\cdot | P_1, \dots, P_K, k_0)$, we have

1003
$$\mathbb{P}(X_{1:\tau} = a_{1:i} | P_1, \dots, P_K) \quad (73)$$

1004

1005
$$= \sum_{k_0 \in [K]} \pi_3(k_0 | P_1, \dots, P_K) \pi_4(a_{1:i} | P_1, \dots, P_K, k_0) \quad (74)$$

1006

1007
$$= \sum_{k_0 \in [K], (P_{k_0})_{1:i} = a_{1:i}} \pi_3(k_0 | P_1, \dots, P_K) \pi_4(a_{1:i} | P_1, \dots, P_K, k_0) \quad (75)$$

1008

1009
$$= \sum_{k_0 \in [K], (P_{k_0})_{1:i} = a_{1:i}} \pi_1(a_{1:i} | P_1, \dots, P_K) m(a_{1:i})^{-1} \quad (76)$$

1010

1011
$$= m(a_{1:i}) \pi_1(a_{1:i} | P_1, \dots, P_K) m(a_{1:i})^{-1} = \pi_1(a_{1:i} | P_1, \dots, P_K). \quad (77)$$

1012

1013 Thus, sampling from π_3 and then π_4 is equivalent to sampling from π_1 , so sampling from π is
1014 equivalent to sampling in the order π_3, π_4, π_2 . Now, let π_5 be the distribution formed by sampling
1015 from π_4 and then π_2 . Then sampling from π is equivalent to sampling in the order π_3, π_5 , where π_3
1016 is a path selection rule conditioned on P_1, \dots, P_K , and π_5 returns $(X_{1:\tau}, Y)$ conditioned only the
1017 selected path and P_1, \dots, P_K . Now, denote q^{π_3} as the true distribution of the path select by π_3 . Then
1018 the output of Φ follows the distribution induced by sampling a *single* path from q^{π_3} , and sampling
1019 from π_5 conditioned only on that single path (the conditioning over P_1, \dots, P_K will cancel). To meet
1020 target-matching and prefix-matching, π_5 must correspond to a valid single-path verification algorithm
1021 with draft distribution q^{π_3} and target distribution p . This completes the proof of the decomposition
1022 into a path selection rule π_3 and a valid single-path verification algorithm with draft values being the
1023 true distribution of the π_3 -selected path. \square
1024

1026 E PROOF OF SKEWED DRAFT DISTRIBUTION FEASIBILITY

1028 Here, we prove [Lemma 4.5](#), which describes when a distribution can be realized as a skewed draft
 1029 distribution q^Γ from a given draft q and path count K . The proof is a direct application of [Lemma A.1](#).
 1030

1031 *Proof.* Define the bipartite graph G with left vertices $(\mathcal{V}^L)^K$, right vertices \mathcal{V}^L , and edges E
 1032 consisting of pairs $((P_1, \dots, P_K), P_k)$ for all $P_1, \dots, P_K \in \mathcal{V}^L$ and $k \in [K]$. Our distribution over
 1033 left vertices is i.i.d. sampling from q , and our distribution over right vertices is r . We would like to
 1034 find when there is a joint distribution on $(\mathcal{V}^L)^K \times \mathcal{V}^L$ supported on E , whose marginals are the left
 1035 and right vertex distributions. [Lemma A.1](#) shows such a distribution exists if and only if
 1036

$$1037 \sum_{P \in S} r(P) \leq \sum_{(P_1, \dots, P_K) \in N(S)} \prod_{k=1}^K q(P_k) \quad \forall S \subseteq \mathcal{V}^L. \quad (78)$$

1040 Here, $N(\cdot)$ denotes neighborhoods in G . Thus, $N(S)$ consists of all (P_1, \dots, P_K) where some
 1041 $P_k \in S$. This means the complement of $N(S)$ is precisely $(\mathcal{V}^L \setminus S)^K$, so the above becomes
 1042

$$1043 \sum_{P \in S} r(P) \leq 1 - \left(\sum_{P \in \mathcal{V}^L \setminus S} q(P) \right)^K \quad \forall S \subseteq \mathcal{V}^L. \quad (79)$$

1046 We are almost at our desired condition, except we need to show this holds for all antichains $S \in$
 1047 $\mathcal{A}(\mathcal{V}^L)$. Indeed, for an antichain S , consider the set $P_S \subseteq \mathcal{V}^L$ of length- L paths with a prefix in the
 1048 antichain. No path in P_S can have two distinct elements in S as a prefix. Thus, the mass of q over S
 1049 is the same as the mass of q over P_S :

$$1050 \sum_{a_{1:i} \in S} q(a_{1:i}) = \sum_{a_{1:i} \in S} \sum_{a_{i+1:L} \in \mathcal{V}^{L-i}} q(a_{1:L}) = \sum_{a_{1:L} \in P_S} q(a_{1:L}) \quad (80)$$

1053 The same holds for r , because it is also a distribution over \mathcal{V}^L . Thus, the antichain conditions reduce
 1054 to the conditions for $S \subseteq \mathcal{V}^L$ above, as desired. \square

1056 F PROOF OF OPTIMAL MULTI-PATH ALGORITHM DESCRIPTION

1058 Now, we prove [Theorem 4.6](#), which describes the optimal multi-path valid verification algorithm as a
 1059 solution to a nonlinear optimization problem.

1061 *Proof.* The first part of the theorem is a direct consequence of the fact that block verification is the
 1062 optimal single-path valid verification algorithm ([Theorem 3.4](#)), and any valid multi-path verification
 1063 algorithm can be decomposed into path selection and valid single-path verification ([Lemma 4.4](#)).
 1064 Now, we prove the second part of the theorem, using this description. Recall from the proof of
 1065 [Lemma 4.5](#) in [Appendix E](#) that a distribution q^Γ can be realized from a given draft q and path count
 1066 K with a path selection rule if and only if

$$1068 q^\Gamma(a_{1:L}) \leq 1 - \left(1 - \sum_{a_{1:L} \in T} q(a_{1:L}) \right)^K \quad \forall T \subseteq \mathcal{V}^L. \quad (81)$$

1070 This covers the feasibility condition, so all that remains is to show the objective function in the
 1071 theorem equals the block efficiency for block verification run on q^Γ . By using the expression for
 1072 $D_{\Phi_{BV}} p$ in the proof of [Theorem 3.4](#) in [Appendix C](#), the objective value in the single-path LP
 1073 ([Theorem 3.3](#)) simplifies to

$$1075 \sum_{a_{1:i} \in \mathcal{V}^L} D_{\Phi_{BV}}(a_{1:i}) p(a_{1:i}) = \sum_{a_{1:i} \in \mathcal{V}^L} w(a_{1:i}) q^\Gamma(a_{1:i}), \quad (82)$$

1077 where the weights w are defined as in [Equation \(3\)](#), but with q replaced by q^Γ . Thus, all that remains
 1078 is to show that

$$1079 w(a_{1:i}) q^\Gamma(a_{1:i}) = \min_{0 \leq k \leq i} p(a_{k+1:i} | a_{1:k}) q^\Gamma(a_{1:k}). \quad (83)$$

1080 The proof is by induction on i . For $i = 0$, this is clear as both sides are one. Now, suppose this
 1081 statement holds for i . To show it is true for $i + 1$, observe that
 1082

$$1083 w(a_{1:i+1})q^\Gamma(a_{1:i+1}) = \min \left\{ 1, \frac{p(a_{i+1}|a_{1:i})}{q^\Gamma(a_{i+1}|a_{1:i})} w(a_{1:i}) \right\} q^\Gamma(a_{1:i+1}) \quad (84)$$

$$1085 = \min \{q^\Gamma(a_{1:i+1}), p(a_{i+1}|a_{1:i})q^\Gamma(a_{1:i})w(a_{1:i})\} \quad (85)$$

$$1087 = \min \left\{ q^\Gamma(a_{1:i+1}), p(a_{i+1}|a_{1:i}) \cdot \min_{0 \leq k \leq i} p(a_{k+1:i+1}|a_{1:k})q^\Gamma(a_{1:k}) \right\} \quad (86)$$

$$1089 = \min \left\{ q^\Gamma(a_{1:i+1}), \min_{0 \leq k \leq i} p(a_{k+1:i+1}|a_{1:k})q^\Gamma(a_{1:k}) \right\} \quad (87)$$

$$1091 = \min_{0 \leq k \leq i+1} p(a_{k+1:i+1}|a_{1:k})q^\Gamma(a_{1:k}). \quad (88)$$

1093 This completes the proof. \square
 1094

1095 G PROOF OF OPTIMAL MULTI-PATH LOWER BOUND

1098 We now prove [Lemma 4.7](#), which lower bounds the objective value in [Theorem 4.6](#).

1100 *Proof.* Denote the minimum value of q^Γ/p as

$$1102 \epsilon = \min_{a_{1:i} \in \mathcal{V}^{\leq L}} \frac{q^\Gamma(a_{1:i})}{p(a_{1:i})}. \quad (89)$$

1104 Then for each $a_{1:i} \in \mathcal{V}^{\leq L}$ and $0 \leq k \leq i$, we have the lower bound

$$1106 p(a_{k+1:i}|a_{1:k})q^\Gamma(a_{1:k}) = \frac{q^\Gamma(a_{1:k})}{p(a_{1:k})} \cdot p(a_{1:i}) \geq \epsilon \cdot p(a_{1:i}). \quad (90)$$

1109 Thus, we have the desired lower bound on the objective value:

$$1111 \sum_{a_{1:i} \in \mathcal{V}^{\leq L}} \min_{0 \leq k \leq i} p(a_{k+1:i}|a_{1:k})q^\Gamma(a_{1:k}) \geq \sum_{a_{1:i} \in \mathcal{V}^{\leq L}} \epsilon \cdot p(a_{1:i}) = (L+1)\epsilon. \quad (91)$$

1113 \square

1115 H GREEDY POLYMATROID CONNECTION

1118 Here, we describe the connection between greedy multi-path valid verification algorithms ([Definition 5.1](#)) and the greedy polymatroid algorithm. We start with [Theorem 4.6](#), but only consider the
 1119 feasibility condition:

$$1122 q^\Gamma(a_{1:L}) \leq 1 - \left(1 - \sum_{a_{1:L} \in T} q(a_{1:L}) \right)^K \quad \forall T \subseteq \mathcal{V}^L. \quad (92)$$

1125 As previously mentioned, the right hand side is a submodular function $\psi(T)$ in T . Note that the
 1126 inequality must also be an equality at $T = \mathcal{V}^L$ to ensure that q^Γ is a probability distribution. Thus, to
 1127 find *some* feasible q^Γ satisfying these submodular constraints, we can use the greedy polymatroid
 1128 algorithm from [Schrijver et al. \(2003\)](#), which aims to maximize the sum of q^Γ terms given the above
 1129 constraints. First, we fix any ordering $\{P_1, \dots, P_M\}$ on \mathcal{V}^L . Then, a feasible solution is

$$1130 q^\Gamma(P_i) = \psi(\{P_1, \dots, P_i\}) - \psi(\{P_1, \dots, P_{i-1}\}). \quad (93)$$

1132 This is precisely what [Equation \(19\)](#) reduces to when the path selection rule Γ is induced by the
 1133 *reversed* global ordering $\{P_M, \dots, P_1\}$. Thus, greedy multi-path valid verification algorithms are
 simply outputs of the greedy polymatroid algorithm for the multi-path LP feasibility conditions.

1134 I OTHER MODEL FAMILIES AND DATASETS

1135
 1136 Here, we expand our original evaluations for OPT-6.7B/350M/125M on GSM8K, HumanEval, and
 1137 MATH500. We broaden our model family coverage to include model pairs Qwen-3 32B/0.6B
 1138 and Llama-3 70B/8B, which allows us to test the efficacy of our approach on much larger target
 1139 models. Due to memory constraints, the Llama-3 experiments were run on a Paperspace machine with
 1140 4xA6000 and an Intel Xeon Gold 5315Y CPU (32 cores, 3.20 GHz). We further expand the Qwen-3
 1141 dataset coverage to include 500 prompts from MGSM and ToolBench, to measure performance
 1142 on diverse, non-academic tasks. Our block efficiency results for $L = 8$ and temperature 1.0 target
 1143 sampling are shown in [Table 3](#), and walltime numbers are in [Table 4](#).

1144 Interestingly, for Qwen-3, the single-path block verification baseline, i.e. $K = 1$, is most effective.
 1145 On GSM8K, block efficiency is 4.32 tokens/ M_p -call at $K = 1$, but decreases to 3.91, 3.57, and 3.43
 1146 for $K = 2, 3, 4$ respectively, while walltime increases from 79.6 to 92.8, 100.5, and 106.1 ms/token.
 1147 Similar trends hold across HumanEval, MATH500, MGSM, and ToolBench. Likewise, for Llama-3,
 1148 $K = 1$ again appears to be the optimal setting. Nevertheless, unlike Qwen-3, there can be meaningful
 1149 gains at higher K . On GSM8K, the walltime decreases by nearly 10% from $K = 3$ to $K = 4$.

1150 1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160				
			$K = 1$	$K = 2$	$K = 3$	$K = 4$	
1151 1152 1153 1154 1155 1156 1157 1158 1159 1160		1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	
1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	GSM8K	4.322	3.910	3.571	3.429
			HumanEval	4.204	3.423	3.331	2.960
			Qwen-3 32B/0.6B	4.342	3.917	3.693	3.411
				4.329	3.973	3.468	3.445
				3.090	2.706	2.791	2.681
1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	1151 1152 1153 1154 1155 1156 1157 1158 1159 1160	Llama-3 70B/8B	4.943	4.430	4.271	4.625
			HumanEval	5.993	4.806	4.473	4.370

1161 Table 3: We compute the block efficiency (tokens per target model call) for decoding with GBV, for
 1162 target-draft model pairs Qwen-3 32B/0.6B and Llama-3 70B/8B on up to 500 prompts from each of
 1163 the datasets GSM8K, HumanEval, and MATH500, MGSM, and ToolBench (only the first two for
 1164 Llama-3). Larger numbers are better. In almost all settings, efficiency decreases monotonically as K
 1165 increases from 1 (block verification baseline) to 4.

1166 1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	1167 1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	1168 1169 1170 1171 1172 1173 1174 1175 1176 1177				1168 1169 1170 1171 1172 1173 1174 1175 1176 1177		
			$K = 1$	$K = 2$	$K = 3$	$K = 4$			
1168 1169 1170 1171 1172 1173 1174 1175 1176 1177		1168 1169 1170 1171 1172 1173 1174 1175 1176 1177		1168 1169 1170 1171 1172 1173 1174 1175 1176 1177		1168 1169 1170 1171 1172 1173 1174 1175 1176 1177		1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	
1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	1168 1169 1170 1171 1172 1173 1174 1175 1176 1177	Qwen-3 32B/0.6B	GSM8K	79.638	92.823	100.547	106.123	
				HumanEval	84.166	106.051	109.215	124.554	
				MATH500	81.517	91.133	97.880	106.575	
				MGSM	80.992	90.312	102.296	105.387	
				ToolBench	113.563	132.477	128.405	133.973	
1176 1177	1176 1177	1176 1177	Llama-3 70B/8B	GSM8K	114.123	133.188	136.272	123.545	
				HumanEval	96.063	124.540	131.889	135.186	

1178 Table 4: We compute the walltimes (ms/token) for decoding with GBV, for target-draft model pairs
 1179 Qwen-3 32B/0.6B and Llama-3 70B/8B on up to 500 prompts from each of the datasets GSM8K,
 1180 HumanEval, and MATH500, MGSM, and ToolBench (only the first two for Llama-3). Smaller
 1181 numbers are better. In almost all settings, walltime increases monotonically as K increases from 1
 1182 (block verification baseline) to 4.

1183 In our Llama-3 70B/8B experiments, we profiled draft tree size and saw that it grows roughly linearly
 1184 with K , at around 9, 16, 24, 29 for $K = 1, 2, 3, 4$. Across the K sweep, the fraction of walltime spent
 1185 in draft and target passes remains nearly constant (around 50% draft and 40% target), and the peak
 1186 target and draft KV-cache footprints remain in a narrow band (target is around 80 MB on GSM8K
 1187 and 110 MB on HumanEval, with the draft being around 40 MB for both).

1188 This shows that increasing K scales FLOPs and KV traffic in proportion to the larger tree, rather than
 1189 changing the draft bottleneck. Even if block efficiency is high, end-to-end walltime degrades with K .
 1190

1191 For Qwen-3 32B/0.6B, we performed the same profiling with similar results. Again, draft tree size
 1192 grows approximately linearly with K . The peak KV-cache footprint also stays in a narrow band as
 1193 K increases, with the target ranging from roughly 60 to 120 MB across datasets and the draft from
 1194 around 30 to 50 MB. Further, draft and target pass times also remain nearly constant relative to the
 1195 walltime. However, compared to Llama-3, the key difference is that 70% of time is spent in drafting
 1196 and only around 20% is in the batched target forward pass, with only small variations across K .
 1197

1198 The target forward pass dominates the cost for Llama-3 significantly more than for Qwen-3. We also
 1199 observe worse walltimes for the same block efficiencies. For example, Qwen-3 on ToolBench sees
 1200 essentially the same walltimes for $K = 1$ and $K = 2$ as Llama-3 on GSM8K, even though Llama-3
 1201 has significantly higher block efficiencies than Qwen-3 (around 1.9 and 1.7 higher tokens/ M_p -call
 1202 for $K = 1$ and $K = 2$). Differences in hardware setups may additionally explain and influence these
 1203 relative target and draft costs.
 1204

1204 J TEMPERATURE ABLATIONS

1205 To measure the impact of temperature on walltime speedups in GBV, we evaluate GBV on Qwen-
 1206 3 32B/0.6B for GSM8K (500 prompts) with temperatures 0.2, 0.4, 0.6, 0.8, also including the
 1207 temperature 1.0 run from [Table 3](#). Our results are shown in [Table 5](#).
 1208

1209 Importantly, we find that GBV performs better at low temperature settings, which directly opposes
 1210 block verification trends. Whereas standard block verification (the $K = 1$ rows) walltimes increase
 1211 from under 80 ms/token to over 90 ms/token from temperatures 0.2 to 1.0, GBV with $K = 3$ sees
 1212 its walltimes decrease from over 100 ms/token at temperature 1.0 to under 86 ms/token for lower
 1213 temperatures (0.2, 0.4, 0.6), with its highest throughput at temperature 0.6. These results also show
 1214 that temperature 1.0 sampling achieves over 10% lower throughput than all other GBV $K = 3$
 1215 settings, so our results in [Section 6](#) actually reflect GBV in its worst-performing setting.
 1216

1217 Even as we noted in [Appendix I](#) and [Table 3](#) that standard BV may outperform GBV for Qwen-3
 1218 32B/0.6B, this finding no longer holds at low temperatures. For example, GBV with $K = 3$ achieves
 1219 a throughput of around 5.6% higher tokens/s than BV when the target temperature is 0.2. Even at
 1220 slightly higher temperatures of 0.4 and 0.6, GBV achieves a higher block efficiency than BV, and
 1221 achieves similar throughput. Thus, even for model pairs where GBV with $K > 1$ does not outperform
 1222 BV at temperature 1.0 sampling, it may still be preferable over BV at low temperatures.
 1223

K	M_p	Temp.	Block Efficiency (tokens/ M_p -call)	Throughput (tokens/s)	Walltime (ms/token)
1	1	0.2	3.867	11.034	90.626
	1	0.4	4.313	12.407	80.598
	1	0.6	4.308	12.368	80.851
	1	0.8	4.430	12.761	78.362
	1	1.0	4.322	12.557	79.638
3	3	0.2	4.219	11.661	85.759
	3	0.4	4.333	11.911	83.954
	3	0.6	4.492	12.149	82.311
	3	0.8	4.161	11.341	88.173
	3	1.0	3.571	9.946	100.547

1237 [Table 5](#): We evaluate how temperature impacts the block efficiency, throughput, and walltime for
 1238 GBV with $K = 1, 3$ when run on Qwen-3 32B/0.6B with GSM8K (500 prompts) and $L = 8$. We use
 1239 temperature 0.2, 0.4, 0.6, 0.8, and 1.0 sampling from the target model. While BV ($K = 1$) performs
 1240 best at high temperatures (0.8, 1.0), GBV ($K = 3$) works better at low temperatures (0.2, 0.4, 0.6).
 1241

1242 K DRAFT LENGTH ABLATIONS

1244 Here, we provide ablation experiments on L , the draft block length. We extend our $L = 8$ results for
 1245 Qwen-3 32B/0.6B on GSM8K with temperature 1.0 target sampling to larger draft lengths: $L = 16$
 1246 and $L = 24$. Our results are shown in [Table 6](#).

1248 As L increases, block efficiency improves only modestly, while throughput degrades rapidly. For
 1249 example, at $K = 1$, the block efficiency increases from 4.32 for $L = 8$ to 5.67 and 5.79 for $L = 16$
 1250 and $L = 24$, but the walltime grows from roughly 80 ms/token to 103 ms/token and 140 ms/token.
 1251 The effect is worse for $K \geq 2$, where moving from $L = 8$ to $L = 24$ more than doubles the walltime.

1252 The draft and target costs also shift substantially: the draft time, as a percentage of the end-to-end
 1253 walltime, grows from about 71% at $L = 8$ to 80% and 86% at $L = 16, 24$, while the target share
 1254 drops to as low as 10%. Coupled with our analysis of the draft bottleneck in [Appendix I](#), this indicates
 1255 longer blocks increase aggregate cost by pushing more work into the draft bottleneck.

L	K	Block Efficiency (tokens/ M_p -call)	Walltime (ms/token)	Draft Cost (%walltime)	Target Cost (%walltime)
8	1	4.322	79.638	71.5	21.3
	2	3.910	92.823	71.3	21.0
	3	3.571	100.547	71.1	21.1
	4	3.429	106.123	70.8	21.3
16	1	5.667	102.970	81.7	13.2
	2	4.016	149.775	81.0	13.3
	3	3.915	155.406	80.5	13.6
	4	3.742	163.200	80.1	13.9
24	1	5.794	140.411	86.0	9.9
	2	4.275	196.937	85.0	10.2
	3	3.945	218.631	84.3	10.8
	4	3.628	240.142	83.4	11.5

1272 Table 6: We evaluate GBV with $1 \leq K \leq 4$ for Qwen-3 32B/0.6B temperature 1.0 sampling on
 1273 GSM8K (500 prompts), over varying draft block lengths $L = 8, 16, 24$. We report block efficiency,
 1274 walltime, and draft and target forward pass runtimes as percentages of end-to-end walltime. We find
 1275 that around $L = 8$ is an ideal length for avoiding more pronounced draft bottlenecks.

1278 L MULTI-PATH COMPARISON

1281 To empirically compare GBV to other multi-path methods, we implement six other speculative
 1282 decoding algorithms in our code: standard speculative decoding (Leviathan et al., 2023; Chen et al.,
 1283 2023), block verification (Sun et al., 2024b), NSS (Miao et al., 2024), SpecInfer (Miao et al., 2024),
 1284 SpecTr (Sun et al., 2023), and traversal verification (Weng et al., 2025). SpecInfer provably improves
 1285 NSS and block verification provably improves standard speculative decoding, so we only benchmark
 1286 GBV against SpecInfer, SpecTr, and block verification.

1287 We run comparison experiments for Llama-3 70B/8B on GSM8K (the same 500 prompts as previous
 1288 experiments) for $L = 8$. We test $K = 2, 3$ as well as temperature 0.2, 1.0 sampling from the target
 1289 model. Our block efficiency, throughput, and walltime numbers are shown in [Table 7](#). Due to poor
 1290 preliminary decoding rates and Llama-3 70B memory footprint, we do not run SpecTr and SpecInfer,
 1291 the worst-performing temperature 1.0 algorithms, at temperature 0.2.

1292 For temperature 1.0 sampling, GBV and traversal verification perform significantly better than SpecTr
 1293 and SpecInfer. For both $K = 2, 3$, GBV achieves over 1.1 token/s higher throughput than SpecTr
 1294 and over 1.6 token/s higher than SpecInfer. While traversal verification does outperform GBV in
 1295 these settings, coming out on top, the incremental gains are less significant: GBV achieves around
 0.6 fewer token/s for $K = 2$ and 0.9 fewer token/s for $K = 3$.

M_p	Temp.	K	Method	Block Efficiency (tokens/ M_p -call)	Throughput (tokens/s)	Walltime (ms/token)
1.0	2	2	GBV	4.430	7.508	133.188
			SpecTr	3.716	6.355	157.365
			SpecInfer	3.259	5.603	178.480
			Traversal	4.760	8.114	123.245
1.0	3	3	GBV	4.271	7.338	136.272
			SpecTr	3.540	6.211	161.002
			SpecInfer	3.239	5.686	175.879
			Traversal	4.779	8.260	121.065
0.2	2	2	GBV	6.730	11.415	87.602
			Traversal	5.650	9.773	102.319
0.2	3	3	GBV	5.663	9.775	102.299
			Traversal	5.763	9.962	100.379

Table 7: We compare GBV, SpecTr, SpecInfer, and traversal verification for Llama-3 70B/8B on GSM8K (500 prompts). We report block efficiency and throughput (higher is better), and walltime (lower is better), for $K = 2, 3$ and temperature 0.2, 1.0 sampling from the target model. While traversal verification outperforms other methods for temperature 1.0, GBV achieves over 1.5 tokens/s higher throughput at temperature 0.2.

On the other hand, at lower temperatures like 0.2, the trend flips: GBV now performs better than traversal verification. While throughput and block efficiency gains are negligible for $K = 3$ (only around 0.1 more tokens/ M_p -call and 0.2 higher token/s), the benefits for $K = 2$ are significant: GBV achieves with over 1.6 higher token/s than traversal verification, and over 1.1 higher tokens/ M_p -call. In fact, GBV at temperature 0.2 with $K = 2$ achieves the highest throughput by far out of all tested settings: no other method for Llama-3 70B/8B achieves over 10 tokens/s throughput. This suggests that GBV outperforms traversal verification for low temperatures and sharper p distributions.