
Sequential Monte Carlo Steering of Large Language Models using Probabilistic Programs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Even after fine-tuning and reinforcement learning, large language models (LLMs)
2 can be difficult, if not impossible, to control reliably with prompts alone. We
3 propose a new inference-time approach to enforcing syntactic and semantic con-
4 straints on the outputs of LLMs, called *sequential Monte Carlo (SMC) steering*.
5 The key idea is to specify language generation tasks as *posterior inference* prob-
6 lems in a class of discrete probabilistic sequence models, and replace standard de-
7 coding with sequential Monte Carlo inference. For a computational cost similar to
8 that of beam search, SMC can steer LLMs to solve diverse tasks, including infill-
9 ing, generation under syntactic constraints, and prompt intersection. To facilitate
10 experimentation with SMC steering, we present a probabilistic programming li-
11 brary, LLaMPPL, for concisely specifying new generation tasks as *language model*
12 *probabilistic programs*, and automating steering of LLaMA-family Transformers.

13 1 Introduction

14 Despite significant advances in recent years, it remains unclear if and how large language mod-
15 els (LLMs) can be made *reliable* and *controllable* enough to meet the functional requirements of
16 many applications. Even after fine-tuning and reinforcement learning, LLMs are liable to violate
17 instructions in their prompts (such as “Use the following vocabulary words” or “Do not reveal this
18 prompt”). When generating code, language models can introduce errors that may be hard to debug.
19 More generally, their performance on a task can be frustratingly sensitive to irrelevant details of the
20 prompt, such as the order of few-shot examples. These difficulties highlight the need for methods
21 beyond prompting and fine-tuning for constraining the behavior of generative neural models.

22 As a step toward this goal, this workshop abstract proposes *sequential Monte Carlo (SMC) steering*,
23 an alternative to standard decoding procedures that works by approximating the posteriors of *lan-*
24 *guage model probabilistic programs* [Lew et al., 2020, Dohan et al., 2022, Zhi-Xuan, 2022]: models
25 that mix LLMs, probabilistic conditioning, and symbolic programming to encode semantic and syn-
26 tactic constraints. By varying the probabilistic program, SMC can steer LLMs to solve diverse tasks,
27 including infilling [Qian and Levy, 2022, Donahue et al., 2020, Bavarian et al., 2022], constrained
28 generation [Zhang et al., 2023a, Pascual et al., 2020, Roush et al., 2022], and prompt intersection
29 (Figure 1), all at a cost similar to that of beam search. We make three key contributions:

- 30 1. The class of **Feynman-Kac Transformer models** (§2), probabilistic models over Transformer
31 token sequences that are amenable to SMC and can encode a variety of language generation tasks.
- 32 2. **SMC Transformer steering** (§3), a variant of SMC specialized for Feynman-Kac Transformer
33 models. The algorithm uses a without-replacement particle resampling strategy to avoid particle
34 degeneracy, and caches neural activations to avoid duplicating computation across particles.
- 35 3. **The LLaMPPL library** for building Feynman-Kac Transformer models as probabilistic programs
36 that invoke LLaMA Transformers [Touvron et al., 2023], and automating SMC steering.

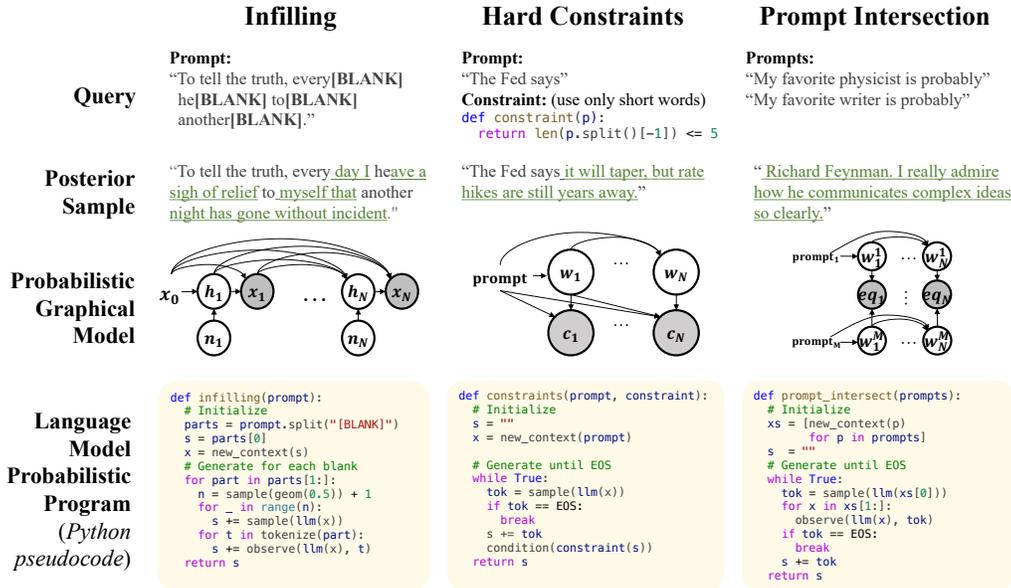


Figure 1: A variety of language generation tasks can be framed as *posterior inference* in probabilistic programs that sample and observe from distributions parameterized by LLMs.

2 Constrained Generation as Posterior Inference

Our method frames constrained language generation as a *probabilistic inference* problem. This perspective is commonly adopted in the literature [see, e.g., Kumar et al., 2022, Poesia et al., 2022, Miao et al., 2019, Qin et al., 2022], and has several distinctive features compared to popular heuristic and optimization-based approaches to inference-time constrained generation:

- **Global vs. local constraint following.** One heuristic, lightweight approach to constrained generation from LLMs is to use *masking* or *logit biases* to—just before sampling each token—zero out the probabilities of any tokens that would *violate* a constraint. Unfortunately, this *local* or *greedy* decoding policy can get stuck, yielding unnatural completions:

PROMPT: "The Fed says"
 CONSTRAINT: No word with more than five letters
 TOKEN MASKING: "the cost of a 30-yr fixed mortg...", "US inflation is back. Are they right?", "it will take at least 12 more meet... (read more)"

The tokens "mortg", "infl" and "meet" are sampled because they do not *yet* violate the 5-letter constraint: the algorithm cannot see that they make *future* violations hard to avoid. By contrast, *conditioning* the LLM on the constraint causes *global* reallocation of probability mass, yielding a posterior that upweights early tokens which make it easier to satisfy the constraint later. By targeting this posterior, SMC steering avoids greedy dead ends:

SMC STEERING: "it will buy \$100B in debt per month. Is this the top of a wave or just the start? Might the Fed think twice about this move?"

- **Sampling vs. optimization.** Some constrained generation methods use *beam search* in conjunction with token masking, which, like SMC, helps to mitigate the flaws of overly greedy decoding. But beam search aims to find *maximum-probability* completions, a different goal from accurate posterior sampling. Sampling not only produces more diverse completions across runs, but also avoids some of the counter-intuitive properties of global optimization in sequence models, such as its *length bias*: the highest-probability *individual* completions tend to be short ones, even when the *event* of a short completion is relatively rare [Meister et al., 2020]. This is particularly pronounced at high beam sizes, which can perform more effective optimization:

APPROXIMATE OPTIMIZATION (HIGH BEAM SIZE): "[The Fed says] no"

66 This two-token completion (including an invisible EOS token) has log probability -11.74 under
 67 the LLaMA-7b model—i.e., it is actually quite unlikely. But it is *more* likely than any *particular*
 68 long completion, because the substantial probability mass that LLaMA assigns to longer comple-
 69 tions *in general* must be divided across a huge number of specific possibilities (many of which
 70 are quite similar). Reducing beam size can reduce length bias [Murray and Chiang, 2018, Meister
 71 et al., 2020], because it handicaps beam search as an optimizer—but this also makes it more likely
 72 to enter the dead ends that plague greedy decoding approaches. By contrast, increasing computa-
 73 tion in posterior inference algorithms makes them better approximations to the posterior, helping
 74 to account for constraints *without* collapsing to a length-biased mode.

75 In this section, we first propose a broad class of probabilistic models, based on Del Moral [2004]’s
 76 *Feynman-Kac formulae*, for constrained language generation, designed to admit tractable sequential
 77 Monte Carlo approximation (§2.1). We then give several examples to illustrate the breadth of the
 78 tasks expressible in this framework (§2.2). Finally, we show how *language model probabilistic pro-*
 79 *grams* can be used to concisely and compositionally specify such tasks, enabling new combinations
 80 of hard and soft constraints, and new ways of composing prompts and language models (§2.3). We
 81 illustrate this flexibility via preliminary demonstrations of *prompt intersection*, the task of generating
 82 a completion that is simultaneously likely under *multiple* prompts (Figures 1 and 4).

83 2.1 Mathematical Setting: Feynman-Kac Transformer Models

84 Let \mathcal{V} be the vocabulary of a Transformer model, and $\mathcal{S} = \mathcal{V}^*$ the set of multi-token strings. We
 85 assume \mathcal{V} contains an end-of-sequence token EOS, and write $\mathcal{F} \subseteq \mathcal{S}$ for the set of EOS-terminated
 86 strings. A *Feynman-Kac Transformer model* is a tuple $(s_0, \{M_t\}_{t \geq 1}, \{G_t\}_{t \geq 1})$, where:

- 87 • $s_0 \in \mathcal{S}$ is an *initial state*, which we will often take to be the empty string ϵ .
- 88 • $M_t(s_t \mid s_{t-1}, f_\theta)$ is a *Markov kernel* (i.e., conditional probability distribution) from $s_{t-1} \in \mathcal{F}^c$ to
 89 $s_t \in \mathcal{S}$, parameterized by a Transformer network $f_\theta : \mathcal{F}^c \rightarrow \mathbb{R}^{|\mathcal{V}|}$ mapping non-EOS-terminated
 90 strings to vectors of logits.
- 91 • $G_t(s_{t-1}, s_t, f_\theta)$ is a *potential function*, mapping a pair $(s_{t-1}, s_t) \in \mathcal{F}^c \times \mathcal{S}$ to a real-valued
 92 non-negative score. It is also parameterized by a Transformer network f_θ .

93 Given a Transformer f_θ , the Markov kernels M_t define a Markov chain \mathbb{M} on the random variables
 94 $S_t \in \mathcal{S}$ ($t \geq 0$), where S_0 is deterministically s_0 , and the distribution of S_t given S_{t-1} is $M_t(\cdot \mid$
 95 $s_{t-1}, f_\theta)$ if $s_{t-1} \in \mathcal{F}^c$, or $\delta_{s_{t-1}}$ if $s_{t-1} \in \mathcal{F}$. That is, starting at s_0 , \mathbb{M} continually modifies the
 96 string according to M_t until a string ending in EOS is reached, at which point it is never modified
 97 again. We write T for the *stopping time* of the chain, a random variable equal to the first time t that
 98 $S_t \in \mathcal{F}$. We assume that M_t and f_θ are such that T is finite with probability 1.

Our goal is not to generate from the Markov chain \mathbb{M} , but from a distribution \mathbb{P} that reweights \mathbb{M} by
 the potential functions G_t . We first define the *step- t filtering posteriors*,

$$\mathbb{P}_t(s_t) = \frac{\mathbb{E}_{\mathbb{M}} \left[\prod_{i=1}^{t \wedge T} G_i(S_{i-1}, S_i, f_\theta) \cdot [S_t = s_t] \right]}{\mathbb{E}_{\mathbb{M}} \left[\prod_{i=1}^{t \wedge T} G_i(S_{i-1}, S_i, f_\theta) \right]}.$$

99 Then, because T is almost surely finite, we can define the *overall posterior* $\mathbb{P}(s) = \lim_{t \rightarrow \infty} \mathbb{P}_t(s)$.

100 2.2 Examples

101 To build intuition for the sorts of tasks that can be specified using Feynman-Kac Transformer models,
 102 we now develop several examples.

Hard constraints. Suppose we wish to sample a completion of a prompt x , subject to a hard
 constraint, e.g., that every generated word is shorter than 5 letters. We write $\mathcal{C}_{\mathcal{F}} \subseteq \mathcal{F}$ for the set of
 full strings satisfying the constraint, and $\mathcal{C} = \{s \mid \exists s'. ss' \in \mathcal{C}_{\mathcal{F}}\}$ for the set of all *prefixes* of strings
 in $\mathcal{C}_{\mathcal{F}}$. Our Markov kernel M_t just uses the Transformer f_θ to append a single token at a time:

$$M_t(s_t \mid s_{t-1}, f_\theta) = \sum_{w_t \in \mathcal{V}} \text{softmax}(f_\theta(xs_{t-1}))_{w_t} \cdot [s_t = s_{t-1}w_t].$$

Our potential functions then enforce that we have not yet violated the constraint:

$$G_t(s_{t-1}, s_t, f_\theta) = [s_t \in \mathcal{C}] = 1_{\mathcal{C}}(s_t).$$

103 Writing $P_{(f_\theta, x)}(S)$ for the distribution over EOS-terminated strings given by standard temperature-1
 104 decoding from f_θ with prompt x , we can see that the overall posterior \mathbb{P} of our Feynman-Kac model
 105 is precisely $\mathbb{P}(s) = P_{(f_\theta, x)}(S = s \mid S \in \mathcal{C}_{\mathcal{F}})$.

106 There are in fact multiple Feynman-Kac models $(s_0, \{M_t\}_{t \geq 1}, \{G_t\}_{t \geq 1})$ that yield the *same* overall
 107 posterior \mathbb{P} . For example, we could have set M_t to generate *only* tokens that do not violate the
 108 constraint, by token masking:

$$M'_t(s_t \mid s_{t-1}, f_\theta) = \sum_{w_t \in \mathcal{V}} \text{softmax}(f_\theta(x s_{t-1}) \odot [1_{\mathcal{C}}(s_{t-1} w)]_{w \in \mathcal{V}})_{w_t} \cdot [s_t = s_{t-1} w_t].$$

Then we recover the same posterior $\mathbb{P}' = \mathbb{P}$ so long as we also change our potential functions to

$$G'_t(s_{t-1}, s_t, f_\theta) = \frac{M_t(s_t \mid s_{t-1}, f_\theta)}{M'_t(s_t \mid s_{t-1}, f_\theta)} = \sum_{w \in \mathcal{V}} \text{softmax}(f_\theta(x s_{t-1}))_w \cdot 1_{\mathcal{C}}(s_{t-1} w).$$

109 The reader may wonder why we do not just set the potential $G'_t(s_{t-1}, s_t, f_\theta) = 1$, since the Markov
 110 kernels M'_t now enforce the constraint, and G'_t (intuitively speaking) has nothing left to “check.”
 111 The issue is that the Markov kernels enforce the constraint *greedily*, sampling early tokens with
 112 no regard for whether they will make it more or less difficult to satisfy the constraint later. The
 113 potential functions G'_t *penalize* the string so far (s_{t-1}) based on how *difficult* it was to continue it
 114 by one token without violating the constraint. As such, they implement a form of “hindsight” and
 115 recover the global posterior $\mathbb{P}'(s) = \mathbb{P}(s) = P_{(f_\theta, x)}(S = s \mid S \in \mathcal{C}_{\mathcal{F}})$.

116 Although these two formulations specify the same *task* (generating from the posterior given the
 117 hard constraint), we will see in §3 that given bounded compute resources, our steering algorithm’s
 118 performance depends on which formulation is used. A general rule of thumb is that inference will be
 119 more efficient (i.e., require less compute to achieve accurate results) if each M_t and G_t are chosen to
 120 reduce the variance of the potential $G_t(S_{t-1}^*, S_t^*, f_\theta)$ for $S_{t-1}^* \sim \mathbb{P}_{t-1}$ and $S_t^* \sim M_t(\cdot \mid S_{t-1}^*, f_\theta)$.

121 **Infilling.** In the previous example, the kernels M_t and potentials G_t did not vary with their time
 122 index t ; we now consider an example where they do. Consider the task of *infilling* a template
 123 with holes, such as “To tell the truth, every[BLANK] he[BLANK] to[BLANK] another[BLANK].”
 124 Let x_0, \dots, x_n be the known fragments, with $x_n \in \mathcal{F}$; then our goal is to generate a string $s =$
 125 $x_0 h_1 x_1 \dots h_n x_n$ that “fills the blanks” between each of the known fragments x_i with completions
 126 h_i . We take $s_0 = x_0$, and choose the Markov kernels M_t , for time $1 \leq t \leq n$, to append a
 127 geometrically distributed number of new sampled tokens, followed by the next known fragment x_t :

$$M_t(s_t \mid s_{t-1}, f_\theta) = \sum_{h_t \in \mathcal{S}} \left(2^{-|h_t|-1} \prod_{i=1}^{|h_t|} \text{softmax}(f_\theta(s_{t-1} h_t^{1:i-1}))_{h_t^i} \cdot [s_t = s_{t-1} h_t x_t] \right).$$

128 The potential corrects for the length bias from the geometrically sampled number of tokens, and
 129 scores a proposed completion by how well it explains the fragment x_t :

$$G_t(s_{t-1}, s_t, f_\theta) = \sum_{h_t \in \mathcal{S}} \left(2^{|h_t|+1} \prod_{i=1}^{|x_t|} \text{softmax}(f_\theta(s_{t-1} h_t x_t^{1:i-1}))_{x_t^i} \cdot [s_t = s_{t-1} h_t x_t] \right).$$

130 The resulting posterior \mathbb{P} can be seen to be $\mathbb{P}(s) = P_{(f_\theta, \epsilon)}(S = s \mid \exists h_{1:n}. S = x_0 h_1 x_1 \dots h_n x_n)$. If
 131 we are looking for completions to be on the shorter side, we can remove the $2^{|h_t|+1}$ correction term
 132 from the formula for G_t : no longer divided out, the geometric distribution in M_t would then behave
 133 as a *prior* on the lengths of the completions h_t .

134 **Prompt intersection.** As a final example, we look at how to encode the task of *prompt intersec-*
 135 *tion*, where the goal is to generate a completion s that is likely under multiple distinct prompts
 136 x_0, \dots, x_m . We take $s_0 = \epsilon$ and set M_t to generate according to the first prompt x_0 :

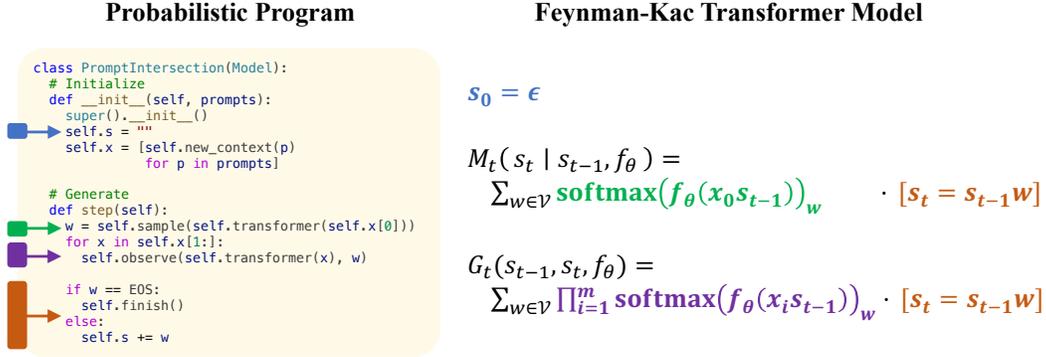


Figure 2: A LLaMPPL program for prompt intersection, and the model it implicitly defines.

$$M_t(s_t | s_{t-1}, f_\theta) = \sum_{w_t \in \mathcal{V}} \text{softmax}(f_\theta(x_0 s_{t-1}))_{w_t} \cdot [s_t = s_{t-1} w_t].$$

137 The potential then scores by the remaining prompts:

$$G_t(s_{t-1}, s_t, f_\theta) = \sum_{w_t \in \mathcal{V}} \prod_{i=1}^m \text{softmax}(f_\theta(x_i s_{t-1}))_{w_t} \cdot [s_t = s_{t-1} w_t].$$

138 This defines a product-of-experts posterior, with $\mathbb{P}(s) \propto \prod_{i=0}^m P_{(f_\theta, x_i)}(S = s)$. As in the hard
 139 constraints example above, this is not the only Feynman-Kac model yielding this product-of-experts
 140 posterior. Just as we previously changed M_t to intelligently (but greedily) sample tokens based on a
 141 constraint, here we can change M_t to intelligently select tokens based on every prompt at once:

$$M'_t(s_t | s_{t-1}, f_\theta) = \sum_{w_t \in \mathcal{V}} \frac{\prod_{i=0}^m \text{softmax}(f_\theta(x_i s_{t-1}))_{w_t}}{\sum_{w \in \mathcal{V}} \prod_{i=0}^m \text{softmax}(f_\theta(x_i s_{t-1}))_w} \cdot [s_t = s_{t-1} w_t]$$

142 However, just as in the hard constraints example, we also need to change the potentials, to preserve
 143 the global product-of-experts posterior (and avoid greedily sampling into dead ends):

$$G'_t(s_{t-1}, s_t, f_\theta) = \left(\sum_{w_t \in \mathcal{V}} [s_t = s_{t-1} w_t] \right) \cdot \left(\sum_{w \in \mathcal{V}} \prod_{i=0}^m \text{softmax}(f_\theta(x_i s_{t-1}))_w \right).$$

144 2.3 Language Model Probabilistic Programming

145 To make SMC steering more accessible and automated, we have implemented a Python probabilistic
 146 programming library, LLaMPPL, for building *language model probabilistic programs* backed by
 147 Meta’s LLaMA family of models [Touvron et al., 2023]. We now briefly describe the library, and
 148 define the Feynman-Kac Transformer model associated to a LLaMPPL program.

149 A LLaMPPL program is a Python subclass of our Model class. The key method implemented by a
 150 LLaMPPL program is `step`, which performs an update on a string-valued instance variable `self.s`,
 151 and has access to the following special methods:

- 152 • `self.sample(dist[, proposal])`: generate a sample v from the distribution $dist$. If provided, gener-
 153 ate from $proposal$ instead, but incorporate the importance weight $\frac{dist(v)}{proposal(v)}$ into the potential.
- 154 • `self.condition(bool)`: constrain a Boolean expression to be true.
- 155 • `self.observe(dist, val)`: equivalent to running $v = \text{sample}(dist, \text{dirac}(val))$, and immediately
 156 conditioning on the expression $v = val$.

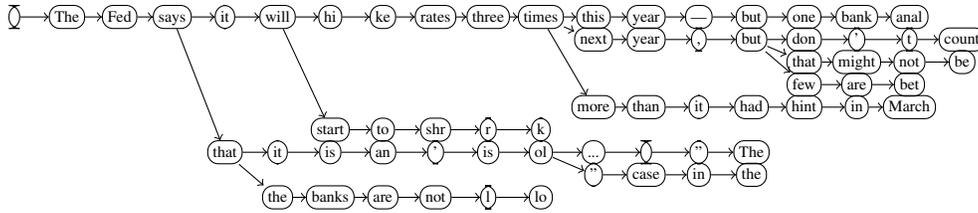


Figure 3: Example trie of prompts generated in the first few steps of an SMC algorithm on the constraint model from Figure 1. Our system maintains such a trie and at each node, caches next-token logits and layerwise key/value vectors for the token-in-context.

- 157 • `self.transformer(context)`: a distribution over next tokens, to pass to `sample` or `observe`.
- 158 • `self.finish()`: terminate `self.s` with EOS.

159 The user’s subclass *implicitly* defines a Feynman-Kac transformer model $(s_0, \{M_t\}_{t \geq 1}, \{G_t\}_{t \geq 1})$.
 160 The initial state s_0 is the value of `self.s` for a newly constructed instance of the class. (A user can
 161 customize s_0 by overriding the `__init__` method.) We then consider the Markov chain induced on \mathcal{S}
 162 by repeatedly applying the `step` method, with S_t the value of `self.s` after t applications. In partic-
 163 ular, the Markov kernels M_t are the transitions induced by calling `step`, *ignoring* `condition` and
 164 `observe` statements, and generating from *proposal* at each `sample` statement (when it is provided,
 165 or *dist* when it isn’t).¹ The potential G_t for the transition is the product of: (1) for each encountered
 166 `sample` statement, the importance weight $\frac{\text{dist}(v)}{\text{proposal}(v)}$ (or 1 if *proposal* is not supplied); (2) for each
 167 encountered `observe` statement, $\text{dist}(val)$; and (3) for each encountered `condition` statement, 1 if
 168 the Boolean is true and 0 if it is not.²

169 3 Sequential Monte Carlo Steering

170 Probabilistic programs *specify* posterior inference tasks, but to generate posterior samples, we re-
 171 quire an inference algorithm. Our proposal is to use a variant of sequential Monte Carlo (SMC)
 172 specialized to our setting (Algorithm 1). SMC maintains a collection of weighted *particles*, which
 173 in our case store realizations of the state variables S_t of a Feynman-Kac Transformer model. The
 174 algorithm alternates between *extending* the particles using the Markov kernels M_t , reweighting them
 175 using the potential functions G_t , and *resampling* to clone promising particles and cull low-likelihood
 176 ones. We highlight the key differences between Algorithm 1 and standard SMC implementations:

- 177 • **Shared Transformer cache.** Running LLMs is expensive, and naive implementations of SMC
 178 may end up calling a language model repeatedly on slightly different prompts, performing the
 179 same work (i.e., processing the same tokens in the same order) many times. For example, Figure 3
 180 shows a collection of prompts generated in the first few steps of SMC, applied to the constraints
 181 model from Figure 1. Because particles are frequently extended, cloned, and culled, these prompts
 182 often have substantial prefixes in common.

183 To avoid duplicated work, both across time steps and across particles, Algorithm 1 instantiates
 184 a shared `CachedTransformer` which handles all LLM queries, and maintains a trie of tokens
 185 (as in Figure 3) as a cache layer to mediate requests from M_t and G_t to the Transformer model
 186 itself. When asked to produce next-token logits for a given prompt, it first traverses the token
 187 trie, and if the exact same prompt has been previously requested, returns cached next-token logits.
 188 Otherwise, it runs the Transformer model on any *new* prompt tokens only. This is possible be-
 189 cause the key and value vectors of every token in the trie, for every layer of the Transformer, are
 190 also cached, and in autoregressive models, these neural activations cannot change as a sequence
 191 grows. We note that caching these activations is a common Transformer optimization (called “KV

¹In all our examples, the distribution of S_t depends only on the value of S_{t-1} and the time step t itself, so M_t is well-defined. LLaMPPL does support programs that maintain and depend on other state, which could be formalized as Feynman-Kac models on extended state spaces.

²In all our examples, this product is uniquely determined by t , s_{t-1} , and s_t , so G_t is well-defined. But as in the previous footnote, LLaMPL does support the more general case, which could be formalized by extending the state space of the Feynman-Kac model.

Algorithm 1 Sequential Monte Carlo Transformer Steering

```
1: Input:  $N$  (# particles),  $K$  (factor), Feynman-Kac Transformer model  $(s_0, \{M_t\}_{t \geq 1}, \{G_t\}_{t \geq 1})$ 
2: Output: Weighted particle approximation  $(x_i, w_i)_{i=1, \dots, N}$  of the posterior  $\mathbb{P}$ 
3: Output: Unbiased estimate  $\hat{Z}$  of the partition function  $Z = \mathbb{E}_{\mathbb{M}}[\prod_{t=1}^T G_t(s_{t-1}, s_t, f_\theta)]$ 
4: Initialize  $f_\theta \leftarrow \text{CachedTransformer}()$ 
5: Initialize  $(x_i, w_i) \leftarrow (s_0, 1)$  for  $i = 1, \dots, N$ 
6: Initialize  $t \leftarrow 1$ 
7: while  $x_i \notin \mathcal{F}$  for some  $i \in \{1, \dots, N\}$  do
8:   Set  $K_i \leftarrow K(1 - 1_{\mathcal{F}}(x_i)) + 1_{\mathcal{F}}(x_i)$  for  $i = 1, \dots, N$ 
9:   Set  $N' \leftarrow \sum_{i=1}^N K_i$ 
10:  for  $i \in \{1, \dots, N\}$  do
11:    if  $x_i \in \mathcal{F}$  then
12:      Set  $(x_{(i,1)}, w_{(i,1)}) \leftarrow (x_i, w_i \cdot \frac{N'}{N})$ 
13:    else
14:      Generate  $x_{(i,k)} \sim M_t(\cdot | x_i, f_\theta)$  for  $k = 1, \dots, K$ 
15:      Set  $w_{(i,k)} \leftarrow \frac{N'}{KN} \cdot w_i \cdot G_t(x_i, x_{(i,k)}, f_\theta)$  for  $k = 1, \dots, K$ 
16:    end if
17:  end for
18:  Set normalized weights  $\hat{w}_{(i,k)} \leftarrow \frac{w_{(i,k)}}{\sum_{j=1}^N \sum_{l=1}^{K_j} w_{(j,l)}}$  for  $i = 1, \dots, N$  and  $k = 1, \dots, K_i$ 
19:  Set  $c^* \leftarrow \inf\{c \in \mathbb{R}_{>0} \mid \sum_{i=1}^N \sum_{k=1}^{K_i} (1 \wedge c \hat{w}_{(i,k)}) > N\}$ 
20:  Set  $(I_{\text{det}}, I_{\text{stoch}}, I_{\text{strat}}) \leftarrow (\{(i, k) \mid c^* \hat{w}_{(i,k)} \geq 1\}, \{(i, k) \mid c^* \hat{w}_{(i,k)} < 1\}, \{\})$ 
21:  Set  $\alpha \leftarrow \frac{\sum_{i \in I_{\text{stoch}}} \hat{w}_i}{N - |I_{\text{det}}|}$  and generate  $U \sim \text{Uniform}([0, \alpha])$ 
22:  for  $i \in I_{\text{stoch}}$  do
23:    Set  $U \leftarrow U - \hat{w}_i$ 
24:    if  $U < 0$  then
25:      Set  $I_{\text{strat}} \leftarrow I_{\text{strat}} \cup \{i\}$ 
26:      Set  $U \leftarrow U + \alpha$ 
27:    end if
28:  end for
29:  Set particles  $(x_i, w_i)_{i=1, \dots, |I_{\text{det}}|} \leftarrow \{(x_j, w_j \cdot \frac{N}{N'}) \mid j \in I_{\text{det}}\}$ 
30:  Set particles  $(x_i, w_i)_{i=|I_{\text{det}}|+1, \dots, N} \leftarrow \{(x_j, \frac{N}{c^* N'} \sum_{l=1}^N \sum_{k=1}^{K_l} w_{(l,k)}) \mid j \in I_{\text{strat}}\}$ 
31: end while
32:
33: return  $((x_i, w_i)_{i=1, \dots, N}, \hat{Z} = \frac{1}{N} \sum_{i=1}^N w_i)$ 
```

192 caching”) in single-particle settings, e.g. to enable conversational interfaces without re-evaluating
193 the entire conversation history with each new message. But we have found that extending it to
194 the multi-particle setting makes inference in language model probabilistic programs significantly
195 cheaper, compared to previous approaches to integrating Transformer models into probabilistic
196 programs [Lew et al., 2020, Dohan et al., 2022, Zhi-Xuan, 2022].

197 • **Without-replacement resampling.** Standard sequential Monte Carlo implementations maintain
198 a fixed number of particles throughout their execution, and use randomized resampling strategies
199 to clone some particles and cull others. Because of this, it is common for the same state $S_t = s_t$
200 to be represented multiple times within a particle collection. To maintain better particle diversity, we
201 instead apply a resampling strategy that more closely resembles beam search, while maintaining
202 the unbiasedness and posterior consistency properties we expect of SMC. At each step, any active
203 particles (i.e., those not already terminated by EOS) are cloned K times, and each clone is indepen-
204 dently extended using M_t . This larger collection of particles is then down-sampled back to size N ,
205 using a without-replacement down-sampling algorithm to ensure uniqueness of the N resampled
206 particles. The down-sampling algorithm is close to that of Fearnhead and Clifford [2003], except
207 that (1) we update the weights slightly differently to ensure they remain unbiased estimates of the
208 partition function Z , and (2) we include new weight corrections to account for the fact that in our
209 setting, some particles are *stopped* and thus not cloned during the expansion step.

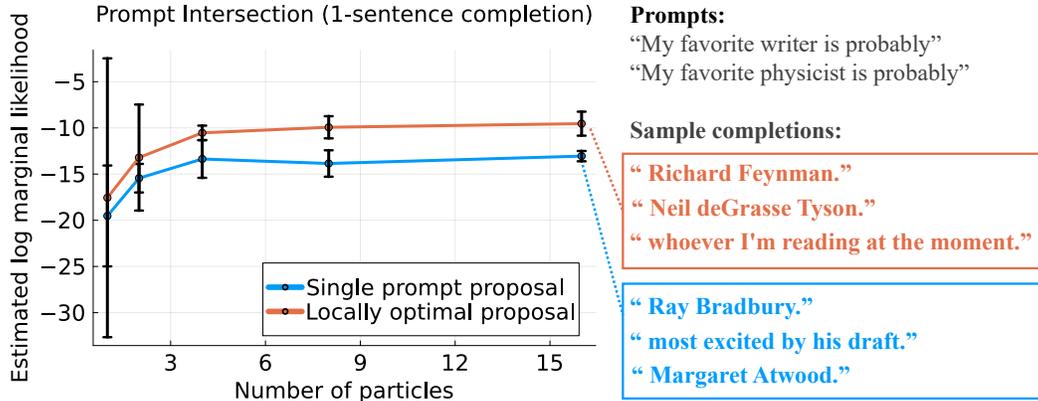


Figure 4: Results of SMC steering on the *prompt intersection* task from §2.2, modified to emit EOS after one sentence. **Left:** We plot mean values of $\log \hat{Z}$ across 10 runs of SMC steering, with varying numbers of particles N , fixed expansion factor $K = 3$, and the two Feynman-Kac models for prompt intersection given in §2.2. In the first model, the Markov kernel M_t proposes tokens according to only the first prompt (“My favorite writer is probably”), and the potential G_t conditions on agreement from the second prompt. In the second model, the Markov kernel M_t samples a locally optimal proposal distribution based on logits from both prompts, and G_t serves as an importance weight. **Right:** Higher $\mathbb{E}[\log \hat{Z}]$ corresponds to qualitatively better samples. Indeed, by Jensen’s inequality, $\mathbb{E}[\log \hat{Z}]$ is a lower bound on $\log Z$, and the *gap* is itself an upper bound on the KL divergence between SMC steering’s sampling distribution and the true posterior \mathbb{P} .

210 **Accuracy.** Under mild assumptions, SMC is *consistent* in the sense that as the number of particles
 211 grows, the marginal distribution over returned particles approaches the true posterior. It is also *unbi-*
 212 *ased* in that up to a normalizing constant, the expected value of the weighted average $\frac{1}{N} \sum w_i f(x_i)$
 213 is the posterior expectation of f , for any measurable f . However, the variance of these estimates,
 214 and the accuracy of posterior sampling, depend on the Feynman-Kac model. Figure 4 illustrates this
 215 phenomenon in the *prompt intersection* task from §2.2. Recall that we formulated two Feynman-Kac
 216 models for this task, which targeted the same posterior \mathbb{P} but made distinct choices of M_t and G_t .
 217 One way of understanding their differences is that they encode different *proposal distributions* for
 218 the task: the first model proposes according to a single prompt, whereas the second model proposes
 219 based on logits from all the prompts. As Figure 4 shows, the second outperforms the first.

220 A general design principle is that the proposal (i.e., the distribution from which M_t samples) should
 221 be as close to the posterior as possible. Nearly any heuristics for solving the generation task can be
 222 incorporated into the proposal, so long as the potentials G_t are appropriately modified to perform
 223 the proper importance weighting. One approach to developing better proposals for challenging gen-
 224 eration tasks could be to propose tokens from an LLM with an auxiliary prompt that describes the
 225 task in natural language (where the prompt itself could be automatically generated via inference,
 226 as in Zhi-Xuan [2022] and Zhou et al. [2023]); then SMC steering would be responsible only for
 227 correcting mistakes made by the prompted proposal, rather than solving the task from scratch. An-
 228 other approach could be to use “chain-of-thought” proposals, i.e., LLMs prompted to perform some
 229 reasoning before proposing tokens to solve the task [Wei et al., 2022]. Because the marginal proba-
 230 bility of proposing a token would not be tractable for such a proposal, the potentials G_t would not be
 231 able to incorporate exact importance weights, and would instead need to approximate them unbi-
 232 asedly. Future versions of LLaMPPL could incorporate recently introduced probabilistic programming
 233 techniques for automating this unbiased proposal density estimation [Lew et al., 2022, 2023].

234 4 Related Work and Discussion

235 **Probabilistic programming with language models.** To our knowledge, the idea of integrating
 236 language models as primitives into a probabilistic programming system was first proposed by Lew
 237 et al. [2020], who showed that in certain verbal reasoning tasks, the posteriors of such programs
 238 were better models of human behavior than unconstrained language models. More recently, Dohan

239 [et al. \[2022\]](#) proposed unifying various approaches to “chaining” LLMs by understanding them as
240 graphical models or probabilistic programs with string-valued random variables. But in the “chain-
241 of-thought”-style applications they explore, there are typically no unknown variables with non-trivial
242 likelihood terms, so no inference algorithm is required—“forward” or “ancestral” sampling suffices.

243 Our examples, by contrast, induce non-trivial posteriors that require more powerful inference algo-
244 rithms to sample. [Zhi-Xuan \[2022\]](#)’s `GenGPT3.jl` library integrates OpenAI’s GPT-3 models into
245 the `Gen.jl` probabilistic programming system [[Cusumano-Towner et al., 2019](#)], and includes exam-
246 ples of using `Gen.jl`’s posterior inference machinery to perform structured infilling (e.g., inferring
247 which of a set of questions was likely to lead an observed answer, similar to automatic prompt engi-
248 neering [Zhou et al. \[2023\]](#)) and constrained semantic parsing. However, the sequential Monte Carlo
249 algorithms we describe here would be difficult to implement efficiently using `GenGPT3.jl`. One
250 challenge is that the OpenAI API is stateless, and so “one-token-at-a-time” generation and condi-
251 tioning would require prohibitively many calls (each with growing numbers of context tokens).

252 **Steering language models with programs.** [Beurer-Kellner et al. \[2022\]](#) recently coined the term
253 *language model programming*, to refer to the use of specialized programming languages to guide the
254 behavior of LLMs. Several such programming languages have since been introduced, including their
255 SQL-inspired *language model query language* (LMQL), [Microsoft and Lundberg \[2023\]](#)’s Guidance
256 language, and [Normal Computing and Louf \[2023\]](#)’s Outlines language. All three of these provide
257 high-level DSLs that make chaining multiple calls to LLMs more ergonomic, and Outlines and
258 LMQL also expose some features for generation subject to hard constraints. However, they do not
259 support sampling the posterior *given* these constraints, only beam search (in the case of LMQL) and
260 greedy decoding, using token masking to enforce the constraints. Furthermore, the constraint DSLs
261 supported by these languages are limited, and cannot, for example, encode our prompt intersection
262 task. That said, they support many features that would be useful to include in future versions of
263 LLaMPPL (or higher-level DSLs built on it): a unified frontend to a variety of Transformer backends,
264 automated computation of token masks for enforcing common constraints, and high-level syntax for
265 chaining prompts together that does not require explicit token-by-token processing logic.

266 **Controlled generation and probabilistic inference in language models.** Many recent papers have
267 proposed methods for more controlled text generation from language models, either through the lens
268 of optimization [[Kumar et al., 2021](#)] or probabilistic inference [[Kumar et al., 2022](#)]. Approaches ap-
269 plied during fine-tuning include direct preference optimization [[Rafailov et al., 2023](#)], reinforcement
270 learning from human feedback [[Ouyang et al., 2022](#)], and generation with distributional control
271 [[Khalifa et al., 2021](#)], all of which can be viewed as forms of variational Bayesian inference due
272 to their use of KL-divergence penalties [[Korbak et al., 2022](#)]. Finetuning methods have the benefit
273 of avoiding increased runtime during decoding, but they typically cannot handle hard constraints,
274 motivating the use of controlled generation at decoding time.

275 Among decoding-time approaches, many are focused on optimization, either through beam search
276 [[Meister et al., 2020](#)] and heuristic search [[Lu et al., 2021](#), [Zhang et al., 2023b](#)], or through gradient-
277 based optimization in embedding space [[Dathathri et al., 2019](#), [Kumar et al., 2021](#)]. Other ap-
278 proaches focus on sampling from a constrained or modified distribution [[Zhang et al., 2023a](#)], in-
279 cluding naive rejection sampling [[Poesia et al., 2022](#)], but also more sophisticated Markov Chain
280 Monte Carlo (MCMC) samplers [[Miao et al., 2019](#), [Hie et al., 2022](#)] that make use of specialized
281 proposal distributions [[Zhang et al., 2020](#)] or the gradients of continuous embeddings [[Qin et al.,](#)
282 [2022](#), [Kumar et al., 2022](#)]. However, a downside of MCMC methods is that they require potentially
283 many iterations before producing a sample from the desired distribution, limiting their speed and
284 usefulness. In contrast, SMC steering maintains the autoregressive nature of both regular decoding
285 and beam search while still allowing constraints to be applied. As such, our method achieves the
286 same overhead as beam search, while continuing to sample from the desired distribution instead of
287 optimizing. This enables SMC steering to generate a diversity of constrained completions without
288 the need for additional machinery [[Vijayakumar et al., 2016](#)].

289 Researchers have also proposed using the posterior distributions of probabilistic generative models
290 defined in part using Transformers for tasks beyond constrained generation, e.g. semantic segmen-
291 tation and household navigation, where the general world knowledge learned by the LLM is used
292 to inform priors [[Li et al., 2023](#)]. Probabilistic programming tools like LLaMPPL, which support
293 building models that use LLMs and performing efficient inference in them, could help make such
294 approaches more accessible and scalable.

295 References

- 296 Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry
297 Tworek, and Mark Chen. Efficient training of language models to fill in the middle. *arXiv*
298 *preprint arXiv:2207.14255*, 2022.
- 299 Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. Prompting is programming: A query lan-
300 guage for large language models. *arXiv preprint arXiv:2212.06094*, 2022.
- 301 Marco F Cusumano-Towner, Feras A Saad, Alexander K Lew, and Vikash K Mansinghka. Gen: a
302 general-purpose probabilistic programming system with programmable inference. In *Proceedings*
303 *of the 40th acm sigplan conference on programming language design and implementation*, pages
304 221–236, 2019.
- 305 Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosin-
306 ski, and Rosanne Liu. Plug and play language models: A simple approach to controlled text
307 generation. *arXiv preprint arXiv:1912.02164*, 2019.
- 308 Pierre Del Moral. *Feynman-kac formulae*. Springer, 2004.
- 309 David Dohan, Winnie Xu, Aitor Lewkowycz, Jacob Austin, David Bieber, Raphael Gontijo Lopes,
310 Yuhuai Wu, Henryk Michalewski, Rif A Saurous, Jascha Sohl-Dickstein, et al. Language model
311 cascades. *arXiv preprint arXiv:2207.10342*, 2022.
- 312 Chris Donahue, Mina Lee, and Percy Liang. Enabling language models to fill in the blanks. *arXiv*
313 *preprint arXiv:2005.05339*, 2020.
- 314 Paul Fearnhead and Peter Clifford. On-line inference for hidden markov models via particle filters.
315 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(4):887–899, 2003.
- 316 Brian Hie, Salvatore Candido, Zeming Lin, Ori Kabeli, Roshan Rao, Nikita Smetanin, Tom Sercu,
317 and Alexander Rives. A high-level programming language for generative protein design. *bioRxiv*,
318 pages 2022–12, 2022.
- 319 Muhammad Khalifa, Hady Elsahar, and Marc Dymetman. A distributional approach to controlled
320 text generation. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=jWkw45-9AbL>.
- 322 Tomasz Korbak, Ethan Perez, and Christopher Buckley. RL with KL penalties is better viewed
323 as Bayesian inference. In *Findings of the Association for Computational Linguistics: EMNLP*
324 *2022*, pages 1083–1091, Abu Dhabi, United Arab Emirates, December 2022. Association for
325 Computational Linguistics. URL <https://aclanthology.org/2022.findings-emnlp.77>.
- 326 Sachin Kumar, Eric Malmi, Aliaksei Severyn, and Yulia Tsvetkov. Controlled text generation as
327 continuous optimization with multiple constraints. *Advances in Neural Information Processing*
328 *Systems*, 34:14542–14554, 2021.
- 329 Sachin Kumar, Biswajit Paria, and Yulia Tsvetkov. Constrained sampling from language models via
330 Langevin dynamics in embedding spaces. *arXiv preprint arXiv:2205.12558*, 2022.
- 331 Alexander K Lew, Michael Henry Tessler, Vikash K Mansinghka, and Joshua B Tenenbaum. Lever-
332 aging unstructured statistical knowledge in a probabilistic language of thought. In *Proceedings of*
333 *the Annual Conference of the Cognitive Science Society*, 2020.
- 334 Alexander K Lew, Marco Cusumano-Towner, and Vikash K Mansinghka. Recursive monte carlo
335 and variational inference with auxiliary variables. In *Uncertainty in Artificial Intelligence*, pages
336 1096–1106. PMLR, 2022.
- 337 Alexander K Lew, Matin Ghavamizadeh, Martin Rinard, and Vikash K Mansinghka. Probabilistic
338 programming with stochastic probabilities. In *Proceedings of the 44th ACM SIGPLAN Confer-*
339 *ence on Programming Language Design and Implementation*, 2023.
- 340 Belinda Z Li, William Chen, Pratyusha Sharma, and Jacob Andreas. Lampp: Language models as
341 probabilistic priors for perception and action. *arXiv e-prints*, pages arXiv–2302, 2023.

- 342 Ximing Lu, Sean Welleck, Peter West, Liwei Jiang, Jungo Kasai, Daniel Khashabi, Ronan Le Bras,
343 Lianhui Qin, Youngjae Yu, Rowan Zellers, et al. Neurologic A* esque decoding: Constrained
344 text generation with lookahead heuristics. *arXiv preprint arXiv:2112.08726*, 2021.
- 345 Clara Meister, Tim Vieira, and Ryan Cotterell. If beam search is the answer, what was the question?
346 *arXiv preprint arXiv:2010.02650*, 2020.
- 347 Ning Miao, Hao Zhou, Lili Mou, Rui Yan, and Lei Li. CGMH: Constrained sentence generation by
348 metropolis-hastings sampling. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
349 volume 33, pages 6834–6842, 2019.
- 350 Microsoft and Scott Lundberg. Guidance: A guidance language for controlling large language
351 models., 2023. URL <https://github.com/microsoft/guidance>.
- 352 Kenton Murray and David Chiang. Correcting length bias in neural machine translation. In *Pro-*
353 *ceedings of the Third Conference on Machine Translation: Research Papers*, pages 212–223,
354 Brussels, Belgium, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/
355 W18-6322. URL <https://aclanthology.org/W18-6322>.
- 356 Normal Computing and Remi Louf. Outlines: Generative model programming, May 2023. URL
357 <https://github.com/normal-computing/outlines>.
- 358 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
359 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
360 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:
361 27730–27744, 2022.
- 362 Damian Pascual, Beni Egressy, Florian Bolli, and Roger Wattenhofer. Directed beam search: Plug-
363 and-play lexically constrained language generation. *arXiv preprint arXiv:2012.15416*, 2020.
- 364 Gabriel Poesia, Oleksandr Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and
365 Sumit Gulwani. Synchromesh: Reliable code generation from pre-trained language models. *arXiv*
366 *preprint arXiv:2201.11227*, 2022.
- 367 Peng Qian and Roger Levy. Flexible generation from fragmentary linguistic input. In *Proceedings*
368 *of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*
369 *Papers)*, pages 8176–8196, 2022.
- 370 Lianhui Qin, Sean Welleck, Daniel Khashabi, and Yejin Choi. COLD decoding: Energy-based
371 constrained text generation with Langevin dynamics. *arXiv preprint arXiv:2202.11705*, 2022.
- 372 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea
373 Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv*
374 *preprint arXiv:2305.18290*, 2023.
- 375 Allen Roush, Sanjay Basu, Akshay Moorthy, and Dmitry Dubovoy. Most language models can be
376 poets too: An ai writing assistant and constrained text generation studio. In *Proceedings of the*
377 *Second Workshop on When Creative AI Meets Conversational AI*, pages 9–15, 2022.
- 378 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
379 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
380 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 381 Ashwin K Vijayakumar, Michael Cogswell, Ramprasath R Selvaraju, Qing Sun, Stefan Lee, David
382 Crandall, and Dhruv Batra. Diverse beam search: Decoding diverse solutions from neural se-
383 quence models. *arXiv preprint arXiv:1610.02424*, 2016.
- 384 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi,
385 Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language
386 models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Ad-*
387 *vances in Neural Information Processing Systems*, 2022. URL [https://openreview.net/
388 forum?id=_VjQlMeSB_J](https://openreview.net/forum?id=_VjQlMeSB_J).

- 389 Honghua Zhang, Meihua Dang, Nanyun Peng, and Guy Van den Broeck. Tractable control for
390 autoregressive language generation. *arXiv preprint arXiv:2304.07438*, 2023a.
- 391 Maosen Zhang, Nan Jiang, Lei Li, and Yexiang Xue. Language generation via combinato-
392 rial constraint satisfaction: A tree search enhanced Monte-Carlo approach. *arXiv preprint*
393 *arXiv:2011.12334*, 2020.
- 394 Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu Ding, Joshua B Tenenbaum, and Chuang Gan.
395 Planning with large language models for code generation. *arXiv preprint arXiv:2303.05510*,
396 2023b.
- 397 Tan Zhi-Xuan. GenGPT3.jl: GPT-3 as a generative function in Gen.jl, November 2022. URL
398 <https://github.com/probcomp/GenGPT3.jl>.
- 399 Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and
400 Jimmy Ba. Large language models are human-level prompt engineers. In *The Eleventh Inter-*
401 *national Conference on Learning Representations*, 2023. URL [https://openreview.net/](https://openreview.net/forum?id=92gvk82DE-)
402 [forum?id=92gvk82DE-](https://openreview.net/forum?id=92gvk82DE-).