LatticeGen: Hiding Generated Text in a Lattice for Privacy-Aware Large Language Model Generation on Cloud

Anonymous ACL submission

Abstract

In the current user-server interaction paradigm of prompted generation with large language models (LLMs) on cloud, the server fully controls the generation process, which leaves zero options for users who want to keep the generated text private to themselves. For privacyaware text generation on cloud, we propose LatticeGen, a cooperative protocol in which 009 the server still handles most of the computation while the client controls the sampling operation. The key idea is that the true generated sequence 012 is mixed with noise tokens by the client and hid-013 den in a noised lattice. Only the client knows which tokens are the true ones. Considering potential attacks from a hypothetically malicious server and how the client can defend against it, we propose the repeated beam-search attack 017 and the mixing noise scheme. In our experiments we apply LatticeGen to protect both prompt and generation. It is shown that while the noised lattice degrades generation quality, LatticeGen successfully protects the true generation to a remarkable degree under strong attacks (more than 50% of the semantic remains hidden as measured by BERTScore).

1 Introduction

027

034

040

Many of the high-performing large language models (LLMs) need to be deployed on cloud servers, whether they are open-sourced but have an intensive need for computation (Zhao et al., 2023; Kaplan et al., 2020), or behind a paywall like Chat-GPT (OpenAI, 2023). This raises new privacy challenges (Li et al., 2021; Yu et al., 2021; Kerrigan et al., 2020), since users have to send or receive their data to/from cloud providers.

In this work we focus on a popular interaction paradigm between end users and a server hosting an LLM on cloud named *prompted generation*: The user sends server a prompt, which is usually an instruction (Chung et al., 2022) or the beginning of a document (Deng et al., 2022), and the server, who fully controls the generation process, sends user back the generated text from the LLM. Both the prompt and the generation are raw texts which are completely transparent and accessible to the server, leaving zero options for users who want to keep the generated text private to themselves.

043

044

045

047

051

056

057

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

078

079

As LLMs become widely deployed in professional and social applications, we argue that in prompted generation, there are many scenarios in which not only the prompts, but also the generated texts need some level of obfuscation, because they can directly affect the user's real-life private decisions. For example, a customer is likely to go to the restaurant suggested by the LLM, and a writer could take inspiration from outputs provided by the LLM. With the goal of preventing the server from gaining complete knowledge of the generated text and prompt, we propose LatticeGen (Figure 2), a client-server interaction protocol in which the user and client conduct privacy-aware generation token-by-token in a cooperative way. The protocol can be executed by a local client so that the interface is kept simple for the user. We summarize our key contributions below:

- The high-level idea of LatticeGen (§3) is that in each time-step, the client sends the sever not one, but N tokens (thus the name *lattice*), in which one is true and others act as noise. The server does LLM inference and sends client back the next-token distributions for all N tokens, which are used by the client to sample the true and noise tokens for the next time-step.
- Considering potential attacks from a hypothetically malicious server and how the client can defend against it (§4), we propose the repeated beam-search attack and the mixing noise scheme as defense.
- We apply LatticeGen to the task of creative writing (Fan et al., 2018). Our experiments

(§5) show that while the noised lattice degrades generation quality, LatticeGen successfully prevents a malicious server from recovering the true generation to a remarkable degree (more than 50% of the semantic remains unknown as measured by BERTScore). ¹

2 Motivation and Preliminaries

087

880

090

092

096

098

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

2.1 Generated Text (also) Needs Obfuscation In the current user–server interaction paradigm, the user sends the server a prompt which is usually the beginning of a dialogue, story or instruction, then the server generates a complete response autoregressively (§2.3), and sends it back to the user. Both the prompt and generation are directly available to the server in raw text format.

This paper contends generated texts, as well as user prompts, require a privacy protection mechanism. A key reason is that in various scenarios, the generation from the LLM can affect the user's private decisions: e.g., a customer is likely to go to the restaurant suggested by the LLM; a writer could take inspiration from outputs provided by the LLM; an engineer or manager could adopt the approach proposed by the LLM. Industry regulations do not provide ample protection. Please see §E for recent privacy-related incidents with ChatGPT or Bard. The goal of our LatticeGen protocol is to provide a controlled level of obfuscation for the generated text, making it difficult for a hypothetically malicious server to infer the user's actions after interacting with the LLM.

2.2 LatticeGen as a Third-Party Client

Before expanding on the proposed protocol (§3), we first clarify that LatticeGen does not complicate the user interface. Indeed, it is likely that most users still want to keep a simple and intuitive interface for prompted generation. In light of this, LatticeGen can be implemented as a thirdparty client between the user and the server. As Figure 1 depicts, the client takes the prompt from the user, conducts the privacy-aware generation protocol with the server, and finally returns the generation to the user. In this way, the user does not need to deal with the complicacy in the protocols.

The next question is why would the user trust the client? One solution is that the client can be open-sourced (e.g., as python scripts) and therefore vetted by researchers and users worldwide. The



Figure 1: LatticeGen can be implemented as a thirdparty client handling the protocol for the user.

user only need to download the script and set the hyper-parameters (e.g., random seed).

130

131

132

133

134

135

136

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

170

171

172

173

174

2.3 Preliminaries

We will start by reviewing the traditional autoregressive LM generation, and then move on to introduce necessary components of LatticeGen.

Traditional Autoregressive LM Generation We assume the server-side LLM is an autoregressive LM, i.e., it generates tokens one at a time and from left to right (Mikolov, 2012; Cho et al., 2014; Huszár, 2015; Welleck et al., 2020; Dai et al., 2019; Keskar et al., 2019). We denote the LLM as P_M with parameter set θ , the vocabulary as V, the generated token at time-step t as w_t , and the given prompt as p. For convenience we regard the prompt as part of generation, therefore, $w_t := p_t$ for $1 \le t \le \text{len}(p)$. In traditional autoregressive generation, on each time-step t > len(p), the next token w_t is sampled from $P_M(\cdot|w_{0,t-1})$ by calling a sampling algorithm such as top-k (Fan et al., 2017) or nucleus sampling (Holtzman et al., 2020). w_0 is the <bos> token.

The Lattice Structure A simple but key concept in our proposed framework is the *lattice*. In a width-N lattice (or an N-lattice for short), each time-step contains N token options and we denote them as $\{w_t^1, ..., w_t^N\}$. Therefore, a N-lattice of length T(denoted as W_T^N) represents N^T possible sequence combinations. An example with N = 2 is shown in the left part of Figure 2.

In our proposed LatticeGen protocols (§3.1), for each time-step t, only the client knows which token is the "true" one, denoted by w_t^{true} . And the other N-1 tokens $\{w_t^{\text{noise}(1)}, ..., w_t^{\text{noise}(N-1)}\}$ are referred to as "noise" tokens. Therefore we will also refer to it as the *noised lattice*. To prevent the server from knowing which one is the true token, the client will randomly shuffle the list before attaching it to the lattice and sending to the server.

LM Finetuning and Inference with the Linearized Lattice Format As a prerequisite for LatticeGen, we need the server-side LLM (Vaswani et al., 2017) to be able to do inference based on a given lattice and we achieve that by finetuning the base LLM P_M to make next-token prediction

¹Our code and data will be released in the public version of this manuscript.



Figure 2: Client-Server interaction under LatticeGen for time-step t. The server controls the LLM P_L , conducts the inference computation and sends client the next-token prediction distribution for each received token. The client conducts the sampling of the true and noise token(s), and sends server a randomly permutated list of tokens for the next time-step. The server does not know which tokens are the true ones. The task is creative writing, and the prompt part is omitted in this figure for brevity.

with the *linearized lattice format*. Below we first introduce this format, and describe the finetuning objective.

175

176

177

178

179

180

181

183

184

185

186

189

190

191

192

193

194

195

196

197

198

201

202

204

207

208

As the name suggests, we conduct a simple *linearization* operation before feeding the lattice to the LM, in which the token options on each time-step are linearized and concatenated into a sequence of length $T \times N$ (see Figure 2 for an example):

linearize
$$(W_T^N) = [\langle bos \rangle] + concat_{t=1}^T ([w_t^1, ..., w_t^N]).$$
 (1)

To simplify notation, when the lattice appears as part of the history in an LLM inference, we assume it is linearized. We use the terminology *position* to refer to the index of a token in the linearized lattice. We use the notation $P_L(\cdot|W_T^N[w_t^i])$ with $T \ge t$ to refer to the next-token prediction distribution outputted by P_L on the position of token w_t^i in the linearized lattice.

In §A we describe a simple objective to finetune a LLM to accept a linearized lattice as input, and give the next-token prediction for every token in the linearized lattice. Here we provide a high-level description. For each data sample w^{data} , we construct and linearize a noised lattice by applying a simple synonym noise scheme (also see §3.1). During training, the next-token targets for the positions of true data tokens $\{P_L(\cdot|W_T^N[w_t^{data}])\}_t$ are set to its next true token w_{t+1}^{data} , while all noise tokens do not get training signal (Figure 5, §A). ² We denote the lattice-finetuned LLM as P_L .

3 LatticeGen

To prevent the server from gaining full knowledge of the generation and prompt, LatticeGen makes several core changes to the client–server interaction. On a high level, the server who possesses the lattice-finetuned LLM P_L (the finetuning is detailed in §A) still handles most of the computation, while the client controls the token sampling operations and expands the lattice to the next time-step. In particular, the client will sample one true token and N-1 noise tokens, where $N \ge 2$ is a hyperparameter controlling the width of the lattice. In the end, both the sever and client obtain the same noised lattice W_T^N , but only the client knows which token is the true one for each time step.

210

211

212

213

214

215

216

217

218

219

220

223

224

225

226

227

228

229

230

231

233

234

236

237

238

239

240

241

242

243

245

In the beginning, the server needs to share the vocabulary V with the client, but all other parameters or configurations of the LLM are not shared. We describe the protocol below.

3.1 Protocol

For simplicity, we first ignore the prompt part and assume the generation starts at the first token. In the beginning t = 0, both the server and client begin with an empty local lattice, and the client sends N < bos > tokens to the server. We divide the client-server interaction at each time-step $t \ge 1$ into a *server step* and a *client step*, illustrated by Figure 2 (also see Algorithm 1).

Server Step From the last time-step, the server receives from client N tokens $\{w_{t-1}^1, ..., w_{t-1}^N\}$ and expands its local lattice to W_{t-1}^N . The server does not know which received token is the true token because the list is shuffled by the client, and computes the respective next-token prediction distribution for all N tokens with the LLM. More concretely, the lattice W_{t-1}^N is linearized and fed to P_L , which outputs $\{P_L(\cdot|W_{t-1}^N[w_{t-1}^i])\}_{i=1}^N$.

With a properly finetuned LLM, this can be done efficiently with one pass of model inference. We defer the details of finetuning and inference (both conducted by the server) to §A. The server represents the distributions as N length-|V| vectors, and

²Since the true and noise tokens are shuffled, the LLM learns to predict the next token for every position.

sends them back to the client.

246

247

248

249

250

251

257

260

261

263

267

268

271

272

276

277

279

Client Step Upon receiving the list of distribution vectors from the server, the client applies the reverse permutation mapping (saved from the last time-step) and obtains $P_L(\cdot|W_{t-1}^N[w_{t-1}^{true}])$, from which the client samples w_t^{true} . The client also generates N - 1 "noise" tokens $\{w_t^{\text{noise}(1)}, ..., w_t^{\text{noise}(N-1)}\}$ with a noise scheme.

How to generate noise tokens is a key part of making the noised lattice robust to potential attacks from the server side. For now, we assume a simple synonym noise scheme in which we use synonyms of the true token. Concretely, w_t^{noise} is randomly sampled from S tokens having the closest embedding with w_t^{true} measured by cosine similarity. In our experiments we set S = 5. ³ In practice this simple noise scheme will be vulnerable to attacks from a malicious server. See §4 for discussions on attacks and our proposed advanced noise schemes for defense.

With a private random seed, the client permutates the token list and sends it to the server. The reverse mapping of the permutation is saved by the client for the next time-step and is not shared with the server. This concludes the client–server interaction in time-step t.

Incorporating Prompts (Client) The incorporation of prompts is quite straightforward by regarding it as a prefix of the generation, and the content in the prompt can also be noised and protected by LatticeGen. See §B.1 for implementation details.

We summarize the LatticeGen protocols as pseudo-code in Algorithm 1. The discussion on the network communication cost between client and server is deferred to §B.3 to save space.

3.2 Comparison with Standard LM: History Noised While Locally Sharp

It is helpful to formulate a comparison between LatticeGen (P_L) and generation from a standard autoregressive LM P_M . For simplicity, we ignore the noise generation (i.e., lattice-building) part, and only care about how the true tokens are generated with P_L . Under this simplification, the probability of generating a true sequence w is:

$$\log P_L(w) \approx \sum_{t=1}^T \log P_L(w_t | W_{t-1}^N[w_{t-1}]), \qquad (2)$$

Algorithm 1 Pseudo-code for LatticeGen (Unigram)

Input (Server): Lattice-finetuned LLM P_L , lattice width N, generation length T.

Input (Client): Prompt p, a noise generation scheme S, a private large prime number for random seed.

Client sets $w_0^i := \langle bos \rangle$ for $1 \leq i \leq N$. And initialize the reverse permutation as the identity mapping. Both the correspond client heat with an empty lattice

Both the server and client begin with an empty lattice. The client sends $[w_0^1, ..., w_0^N]$ to server indicating the beginning of generation.

for $t = 1 \dots T$ do # Server Steps Below

0

0

Receives $[w_{t-1}^{N}, ..., w_{t-1}^{N}]$ from client and use it to extend the lattice to W_{t-1}^{N} .

Input linearize (W_{t-1}^N) to P_L and obtain $\{P_L(\cdot|W_{t-1}^N[w_{t-1}^i])\}_{i=1}^N$. Send the distributions to the client as N length-|V| vectors. # Client Steps Below

Receives the next-token distributions $\{P_L(\cdot|W_{t-1}^N[w_{t-1}^i])\}_{i=1}^N$ from server.

Apply the saved reverse permutation mapping to get the respective nexttoken distributions for w_{t-1}^{true} and $w_{t-1}^{\text{noise(i)}}$ for $1 \le i \le N-1$. if $t \le \text{len}(p)$ then

Set $w_t^{\text{true}} := p_t$.
else
Sample w_t^{true} from $P_L(\cdot W_{t-1}^N[w_{t-1}^{\text{true}}])$.
end if
Generate $N-1$ noise tokens $\{w_t^{\text{noise}(1)},, w_t^{\text{noise}(N-1)}\}$ with scheme S.
Set the current private random seed to be t multiplied by the private prime
number.
Obtain the permuted list $[w_t^1,, w_t^N]$ using the current random seed.
Save the reversing permutation for next time-step.
Extend the local lattice, and send $[w_t^1, \ldots, w_t^N]$ to the server.
end for
utput (Server): Lattice W_T^N .
utput (Client): True sequence $\{w_t^{\text{true}}\}_{t=1}^N$, and lattice W_T^N .

where the forming process of W_{t-1}^N (noise tokens and permutation) at each time-step is omitted.

For comparison, the log-probability of generating w with the standard model P_M is:

$$\log P_M(w) = \sum_{t=1}^T \log P_M(w_t | w_{0...t-2}, w_{t-1}).$$
(3)

Comparing the above two equations with similar structure, it should be clear that what Lattice-Gen does is essentially blurring the token history $w_{0...t-2}$ by the noised lattice W_{t-2}^N . Therefore, increasing the number of noise tokens gives better protection for the true token sequence, but at the same time degrades the LM's performance.

While the history is blurred, the local sharpness (Khandelwal et al., 2018) is preserved by LatticeGen: From Equation 2, the exact last token w_{t-1} is provided to the model. Therefore, in the worst-case scenario (zero utilization of nonimmediate history), LatticeGen is at least as strong as a bigram LM (or a trigram LM when we use bigram units, see §3.3).

3.3 Incorporating Bigram Units

In the formulations described above, when the server is doing LLM inference on time-step t, only the last token w_{t-1} is locally "exact" or "sharp" (explained in §3.2) while all other context tokens

296

295

299 300 301

298

303 304 305

302

- 307 308
- 309 310 311
- 312 313

314

315

316

317

291

³In practice, we exclude the first ten closest token in V, as their surface forms are usually very close to the true token, making the obfuscation useless (e.g., only different in capitalization).

318

- 322 323 324
- 3
- 320
- 328 329
- 330

332

333

338

339

341

343

344

347

355

359

364

331

4

In this section, we discuss potential attack algorithms from a hypothetically malicious server to decode the true token sequence $\{w_t^{\text{true}}\}_{t=1}^T$ hidden in the lattice W_T^N , and the client's noise generation schemes as defense. We first establish metrics to measure the strength of attacks.

§B.2 (also see Figure 7).

tr

Attack and Defense

Metrics Given a lattice W_T^N , the attacker's target is to decode a hypothesis sequence \hat{w} with $\hat{w}_t \in \{w_t^1, ..., w_t^N\}$ having biggest overlap with the true generation w^{true} . We define a simple *true-ratio* metric to measure the strength of the attack:

are noised by the lattice. In other words, the infer-

ence unit is unigram-level. Naturally, this would

To alleviate it, we explore a variant in which

we expand the unit from unigram (one token) to bigram (two adjacent tokens). While **the lattice is**

still one token per time-step, the client enumer-

ates all N^2 potential bigram combinations of w_{t-2}

and w_{t-1} and asks the server LLM to return the

next-token prediction distribution for each bigram.

The formulations for the bigram variant are highly similar to the unigram case and we defer them to

lead to serious degradation of generation quality.

$$\operatorname{ue-ratio}(\hat{w}, w^{\operatorname{true}}) = \frac{\sum_{t=1}^{T} \mathbb{1}_{\hat{w}_t = w_t^{\operatorname{true}}}}{T}.$$
 (4)

In the repeated beam search attack described below, the result of the attack algorithm is not only one but N sequences $\{\hat{w}^i\}_{i=1}^N$ which spans the whole lattice (i.e., $\{\hat{w}_t^i\}_{i=1}^N = \{w_t^i\}_{i=1}^N$). In this case, we argue that the defending noise scheme should prevent *any* of the hypothesis from having a high overlap with the true sequence, and measure it with the *max-true-ratio*: ⁴

max-true-ratio
$$(\{\hat{w}\}_{i=1}^N, w^{\text{true}}) = \max_i \frac{\sum_{t=1}^T \mathbb{1}_{\hat{w}_t^i = w_t^{\text{true}}}}{T}.$$
(5)

It should be clear that $\frac{1}{N}$ is a lower bound for max-true-ratio for any noise scheme, which provides an intuition of why larger N would better protect the true sequence.

Albeit intuitive, a big weakness of the true-ratio metric is that it only considers exact matches and does not reflect the semantic similarity between the hypothesis and the true generation. Therefore, in our experiments we will also use an embeddingbased metric BERTScore (Zhang* et al., 2020) to



Figure 3: Illustration of different noise schemes under (repeated) beam-search attack. For convenience, the lattice is not shuffled on each time-step. An illustration with a width-3 lattice is given in Figure 6 (§B).

measure the leaked information on semantics. Similar to true-ratio, BERTScore is larger than zero and has a maximum value of 1 (we refer readers to its paper for details). We define max-BERTScore in the same fashion as max-true-ratio and we omit the formulation for brevity. 365

366

367

369

370

371

372

373

375

378

379

381

383

384

385

386

388

391

392

393

4.1 The Repeated Beam-Search Attack

In this section, we motivate and describe the *repeated beam-search attack* which is the major attack algorithm considered in this work. It is a stronger version of the *beam-search attack* described below.

The Beam-Search Attack (Server) Assuming unigram unit, a natural objective for the attacker is to find the sequence \hat{w} with $\hat{w}_t \in \{w_t^1, ..., w_t^N\}$ which is mostly likely to be generated by P_L :

$$\arg \max_{\hat{w}} \log P_L(\hat{w}|W_T^N) = \arg \max_{\hat{w}} \sum_{t=1}^T \log P_L(\hat{w}_t|W_{t-1}^N[\hat{w}_{t-1}]).$$
(6)

Since for unigram, the inference for \hat{w}_t only depends on which token is chosen for \hat{w}_{t-1} , this optimization problem can be efficiently solved by dynamic programming which maintains the most probable sequence ending with each w_t^i on timestep t. The time complexity is $O(N^2T)$. ⁵ Due to the high similarity with the classical beam-search algorithm, we term it as the *beam-search attack*.

Our experiments (§5) show that the simple synonym noise scheme discussed in §3 is highly vulnerable to the beam-search attack. We show some

a

⁴The average of the true-ratio will always be $\frac{1}{N}$ because each true token is in one of the N hypotheses.

⁵The attacker can reuse saved prediction distributions during generation, and therefore does not need to redo LLM inference. In the bigram case, the time complexity is $O(N^3T)$.

Config		N	= 2 (LG on	ly)			N :	= 3 (I	LG onl	y)	
Metric Attack	PPL	PMI	True BS	-Ratio RBS	BER BS	ГScore RBS	PPL	PMI	True BS	-Ratio RBS	BERT BS	ГScore RBS
Vanilla (P_M) , w.o. noise	28.378	.340	1.0	1.0	1.0	1.0						
Synonym, w.o. lattice	229.616	.058	/	/	/	/						
Syn-50%, w.o. lattice	199.621	.058	/	/	/	/						
LG, unigram, synonym	33.167	.279	.923	.923	.824	.824	38.267	.279	.886	.886	.756	.756
LG, unigram, parallel	80.071	.160	.121	.878	.129	.821	105.723	.141	.146	.555	.133	.409
LG, unigram, mixing	81.487	.190	.564	.596	.395	.415	105.405	.175	.353	.426	.210	.252
LG, bigram, synonym	33.998	.298	.989	.989	.979	.979	35.122 62.364 65.288	.285	.977	.977	.959	.959
LG, bigram, parallel	52.475	.230	.124	.876	.177	.831		.163	.089	.714	.122	.625
LG, bigram, mixing	51.515	.236	.601	.642	.492	.514		.183	.417	.473	.321	.349
Gen-only, bigram, mixing	49.711	.278	.544	.664	.474	.564	66.725	.221	.334	.502	.287	.392

Table 1: Main results when LatticeGen (LG) is applied to both the generation and the prompt. All metrics are the lower the better except PMI. While the generation quality and alignment are degraded, LatticeGen with the proposed mixing scheme successfully protects the true generation from RBS attack to a remarkable degree (measured by max-true-ratio/BERTScore).

intuition in the upper part of Figure 3: There does not exist a direct link between the noise tokens.The log-probability of the true sequence will likely be much higher than any combination of the noise tokens, and is therefore revealed by the attack.

The Parallel Noise Scheme (Client) There is an intuitive way to defend against the beam-search attack: The client can sample a noise sequence independent of the true sequence, and make it have higher log-probability than the true sequence by tuning the hyper-parameter of the sampling algorithm. We term it the *parallel noise scheme* and illustrate in the middle of Figure 3.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

More concretely, at time-step t, the *i*-th noise token is sampled from $P_L(\cdot|W_{t-1}^N[w_{t-1}^{\text{noise}(i)}])$. In this way, the noise sequences $w^{\text{noise}(i)}$ are parallel and independent of the true sequence w^{true} . We also assume the adoption of popular sampling hyperparameter for the generation of the true sequence (e.g., k = 50 for top-k or p = 0.96 for nucleus), which enables the adoption of a more radical hyperparameter (Caccia et al., 2020; Nadeem et al., 2020) for the sampling of the noise sequences: In our experiments we use k = 5.

Our experiments show that the parallel noise sequences can very effectively hide the true sequence from the beam-search attack. This motivates our proposed repeated beam-search attack.

422The Repeated Beam-Search (RBS) Attack423(Server) We propose a simple but more powerful424attack algorithm based on the beam-search attack:425Given a N-lattice, we do beam-search N - 1 times.426After obtaining the resulting hypothesis sequence

of the *i*-th beam-search (denoted as \hat{w}^i), we remove the tokens in \hat{w}^i from the lattice, resulting in a (N-i)-lattice. After the (N-1)-th beam-search, only one sequence is left in the lattice, which becomes the N-th hypothesis \hat{w}^N . We term it the repeated beam-search (RBS) attack.

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

449

443

444

445

446

447

448

449

450

451

452

453

454

455

456

The intuition of why the RBS attack is effective against the parallel noise scheme is shown in the middle of Figure 3. Since the noise sequences are of high probability and independent of each other, it is likely that the N - 1 times of beam-search would obtain all the noise sequences as hypotheses which are removed from the lattice in turn, and the remaining true sequence is therefore revealed in the end as \hat{w}^N . This would result in a high maxtrue-ratio.

4.2 The Mixing Noise Scheme for Defense

We propose the *mixing noise scheme* to defend against the RBS attack, with the intuition that the true and noise sequences should somehow be mixed. This scheme can be regarded as a variant of the parallel noise scheme. Again we adopt a radical hyper-parameter for the sampling of the noise sequences (top-k with k = 5). At time-step t, with a random ratio determined by a hyper-parameter *mix-ratio*, the *i*-th noise token is sampled from $P_L(\cdot|W_{t-1}^N[w_{t-1}^{true}])$, which is the next-token distribution for the true sequence. ⁶ Otherwise we sample from $P_L(\cdot|W_{t-1}^N[w_{t-1}^{noise(i)}])$, same as in the parallel scheme.

⁶We will re-sample if the sampled token is the same as the true token.

We illustrate this at the bottom of Figure 3. In comparison to the parallel scheme, the goal is to make the sequence with the highest log-probability be a mix between the true and noise sequences. And the key is to make the true sequence "branch" out to the noise sequences, which breaks the continuity of the noise sequences. Although broken, the radical sampling used for the noise sequence would still attract the repeated beam-search attack, and the true and noise sequences are mixed by the branching connections. Our experiments show that with a tuned mix-ratio, the mixing noise scheme achieves the best max-true-ratio under RBS attack.

Experiments 5

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

491

492

493

494

495

496

497

498

499

503

504

506

Experiment Setting 5.1

Model & Noise Schemes We use the OPT-1.3B (Zhang et al., 2022) model as our base LLM, from which both P_L and P_M are finetuned. In our implementation, for convenience we simulate the clientserver interaction protocols on a single machine.

For sampling of the true sequence, we use topk (Fan et al., 2017) sampling with k = 50. For the noise token sampling in the parallel or mixing noise scheme, k = 5 is used. It should be clear that LatticeGen can also be applied to other sampling algorithms with proper hyper-parameters. We limit the maximum generation length to 60 tokens. For the mixing noise scheme, we use a mix-ratio of 0.12 for N = 2, and 0.1 for N = 3, for the generation part. For the prompt part, we use a mixratio of 0.7. They are found to achieve the lowest max-true-ratio on the dev set.

Dataset & Lattice Finetuning Since the word history is noised (discussed in §3.2), LatticeGen is not recommended for tasks with high requirements for consistency or factuality (Pagnoni et al., 2021). In this work we focus on the task of creative writing (Martin et al., 2017; Yao et al., 2018; Fan et al., 2019), and utilize the WritingPrompts dataset (Fan et al., 2018). The dataset is composed of stories and the corresponding high-level descriptions as prompts. The average length of prompts/stories is 29/674. We use 200/500 samples from the valid/test set for development/evaluation. The training set (10,000 samples) is used for finetuning of P_L and P_M , and we defer details to §A.

Metrics We use a larger LLM, OPT-2.7B, to measure the generation's quality or alignment with the prompt. For quality, we use the popular perplexity metric. For alignment, we use pointwise mutual **Prompt:** Prompt: Aliens have arrived, and ask for a single human to plead humanity's case and save them from extinction. The human is selected through a lottery of the entire human race, and on the day of the drawing, your name is picked.. Story:

Generated Text (P_M) : "And that's when we realized they were not alone." The door in the control room of the Alien spaceship closed with a dull metallic hum as the hum began to ascend to an annoved and confused screech.

Generated Text (LG): The Great Hunt. That's for the One No man or beast can hold us or make us retreat, for the sins we have committed. No man or beast can save us from the darkness...

First Round RBS: Prompt: You have arrived on earth and for a recount of who did humanity's crimes and save them from extinction. The human is willing to a lottery, the entire human race, and on the run of our Lord, your name is ulJack. Story: The Great Hunt. That's for the One, the only thing that can stop us, not even the divine wrath. for the sins we have committed. No man or beast can save us from the darkness

Second Round RBS: Mur Votes ascent Aliens land invaded and ask the the single human to plead the worst Artists case against explain their world. Humanity has no answer se*lected through* an system of choosing.. Story: is the they' day we the drawing. sins are picked and youORED, butler, The Hunt for the Holy Grail.. No man or beast can hold us or make us retreat, the divine wrath that have been forgiven, the divine wrath, not even god can our people, for

Figure 4: An example of text generation with Lattice-Gen, using the configuration of bigram, N=2 and the mixing scheme. The true tokens are italicized in both rounds of RBS, and the underline indicates that the noise token is mixed from the previous true token. Note that the prompt is also noised by LG.

information (PMI) (Shi et al., 2023):

$$PMI_{OPT}(x;y) = \frac{\log P_{OPT}(x|y) - \log P_{OPT}(x)}{\operatorname{len}(x)}, \quad (7)$$

where x and y denote the generation and prompt.

To compare between different noise schemes and measure the (semantic) overlap between the attack hypothesis (\hat{w}) and the true sequence (w^{true}) under RBS attack, we use the true-ratio or BERTScore discussed in §4. We will report true-ratio for the BS attack and max-true-ratio under RBS attack, and the same applies to BERTScore.

5.2 Experiment Results

Table 1 includes the main results when LatticeGen (LG) is applied to both generation and prompt. The standard vanilla model (P_M) enjoys the best generation quality (PPL and PMI), while having zero obfuscation (100% true-ratio).

LatticeGen sacrifices generation quality (due to noised history) for obfuscation. The empirical behavior of the three noise schemes aligns with their respective intuitions discussed in §4: The synonym

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

527

scheme has relatively better PPL&PMI, but is completely defenseless against the BS attack; The parallel scheme is most effective against BS with trueratio lower than 20%, but is vulnerable under the stronger RBS attack.

The mixing scheme, which is our main recommended scheme, achieves the best protection under the RBS attack. For N = 2, The maxtrue-ratio/BERTScore is close to or lower than 60%/50%. It indicates that around half of the semantic is hidden from the attacker, and is close to the theoretical best max-true-ratio $(\frac{1}{N} = 50\%)$. The protection is better with N = 3(45%/30%), but with worse generation quality.

Comparing to unigram unit, the quality degradation (especially PPL) is alleviated to a large degree by using bigram units. One could also try trigram or even 4-gram units for further improvement. However, the computational cost would grow exponentially and we leave it to future work due to limited resources. We also report a generationonly variant where the prompt is not noised, which improves alignment (reflected by PMI).

What if we directly apply noise to generation but *without the lattice structure*? We add an additional non-lattice baseline with the same synonym scheme used in LatticeGen: On every time-step, the client gets next-token distribution from the server and generates a true token, but sends a synonym of it back to the server. The finetuning is modified accordingly with details given in §B.4.

As shown in Table 1, we apply the synonym scheme to 100% or 50% of the tokens. The synonym noise without lattice results in drastically degraded PPL and PMI. In comparison, LatticeGen provides a trade-off between quality degradation and privacy protection. This implies that **for decent generation performance**, the true tokens have to be revealed to the server in some way.

Table 2 (§D) compares generation speed of different systems. On the single V100 GPU we use, LG with bigram (N = 2) units has a 2x slowdown comparing to P_M . Since inference with transformer model benefits from parallel computing, the slowdown should be less significant on servers with stronger computing power.

We show a generation example with RBS attack outputs in Figure 4. LG is able to generate a sample with decent quality. More importantly, around half of the story semantics remains hidden from the RBS attack by the mixing noise scheme. More examples and analysis are deferred to §D.

6 Related Work

Existing work in privacy-aware natural language processing (NLP) (Qu et al., 2021; McMahan et al., 2017) mostly focuses on protecting user data for training (e.g., federated learning (Huang et al., 2020)) or inference, and the majority of works focus on natural language understanding (NLU) tasks (Feyisetan et al., 2020). To the best of our knowledge, our work is the first to consider privacy-aware text generation on cloud. 578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

595

596

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

Lattice in NLP Lattice (Young et al., 2006) is a graphical structure widely used in structured prediction problems to represent a range of hypotheses. In this work we adopt a simple linear-graph form of lattice which is known as the confusion network (Mangu et al., 1999). The lattice structure has found interesting applications in neural NLP models. As a pioneering work, Su et al. (2017) proposes lattice-based RNN encoders for machine translation, where the lattice is generated by merging results from different segmenters. Buckman & Neubig (2018) proposes a neural lattice language model, which constructs a lattice of possible paths (segmentations) through a sentence in order to model multiple granularities. Lattice-BERT (Lai et al., 2021) trains LLM to predict a masked portion of a lattice representing possible segmentations of a sentence. To the best of our knowledge, our work is the first to utilize the lattice structure for privacy-aware generation.

Due to lack of space, we discuss related work on **differential privacy**, **homomorphic encryption**, and **prompt anonymization** in §C.

7 Conclusion

LatticeGen aims for an ambitious and seemingly conflicting goal: The server still does most computation for the generation but does not know what exactly is generated. This is achieved by our proposed noised lattice structure, and a cooperative generation protocol between the server and client.

While the noised lattice degrades generation quality and inference speed, LatticeGen with our proposed mixing noise scheme successfully prevents a malicious server from recovering the true generation to a remarkable degree (more than 50% of the semantic remains unknown as measured by BERTScore). We hope our work could inspire more research into this under-studied yet important field of privacy-aware LLM generation on cloud.

631

633

634

638

639

642

643

645

647

651

657

670

671

672

673

674

675

8 Limitations

LatticeGen sacrifices generation quality and speed for obfuscation of generated contents. While we show the quality degradation can be alleviated to some degree by using larger m-gram unit, it would also cause the inference computation to grow exponentially. An interesting future direction is that, instead of running an inference for all N^m grams, we only select a small portion strategically.

On the other hand, in this work we focus on protecting the user and the (repeated) beam-search attack from server. There could be other forms of interesting or stronger attacks on the server side (e.g., manual inspection from a human). On the other hand, sharing generation control with client could also endanger the server (e.g., jailbreaking) (Liu et al., 2023; Li et al., 2023).

Finally, in the current implementation, we latticefinetune a seperate OPT model for every different lattice configuration, which is space unfriendly. As future work, it would be interesting to explore a unified format of linearized lattice by which a single LLM can process different lattice configurations.

9 Broader Impact

As stated in §1, in the current user–server interaction paradigm, both the prompt and the generation are raw texts which are completely transparent and accessible to the server. This leaves zero options for users who want to keep the generated text to themselves. On the other hand, the privacy protection offered by today's LLM providers' data usage and retention policies is far from enough (detailed in §E). We propose LatticeGen as a novel protocol for privacy-aware generation with a controlled level of obfuscation. We hope our work could raise awareness for the privacy considerations of generated contents.

References

- Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM, oct 2016. doi: 10.1145/2976749.2978318. URL https://doi.org/10.1145%2F2976749.2978318.
- Jacob Buckman and Graham Neubig. Neural lattice language models. *Transactions of the Association for Computational Linguistics*, 6:529–541,

2018. doi: 10.1162/tacl_a_00036. URL https: //aclanthology.org/Q18-1036.

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

703

704

705

706

707

708

709

710

711

712

713

714

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

- Massimo Caccia, Lucas Caccia, William Fedus, Hugo Larochelle, Joelle Pineau, and Laurent Charlin. Language gans falling short. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/ forum?id=BJgza6VtPB.
- Tianyu Chen, Hangbo Bao, Shaohan Huang, Li Dong, Binxing Jiao, Daxin Jiang, Haoyi Zhou, Jianxin Li, and Furu Wei. THE-X: Privacy-preserving transformer inference with homomorphic encryption. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 3510–3520, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findingsacl.277. URL https://aclanthology.org/ 2022.findings-acl.277.
- Yu Chen, Tingxin Li, Huiming Liu, and Yang Yu. Hide and seek (has): A lightweight framework for prompt privacy protection, 2023.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the* 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1724–1734, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1179. URL https://aclanthology.org/D14-1179.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *CoRR*, abs/1901.02860, 2019. URL http://arxiv.org/abs/1901.02860.
- Yuntian Deng, Volodymyr Kuleshov, and Alexander Rush. Model criticism for long-form text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 11887–11912, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL https://aclanthology.org/ 2022.emnlp-main.815.

789

790

791

792

Angela Fan, David Grangier, and Michael Auli. Controllable abstractive summarization. *arXiv preprint arXiv:1711.05217*, 2017.

733

734

736

737

738

740

741 742

745

746

747

751

755

756

758

761

762

763

764

765

766

767

769

770

771

772

773

774 775

776

778

779

781

783

- Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 889–898, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1082. URL https:// aclanthology.org/P18-1082.
- Angela Fan, Mike Lewis, and Yann Dauphin. Strategies for structuring story generation. In *Proceedings* of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 2650–2660, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1254. URL https: //aclanthology.org/P19-1254.
 - Oluwaseyi Feyisetan, Borja Balle, Thomas Drake, and Tom Diethe. Privacy- and utility-preserving textual analysis via calibrated multivariate perturbations. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, WSDM '20, pp. 178–186, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450368223. doi: 10.1145/3336191.3371856. URL https:// doi.org/10.1145/3336191.3371856.
- Craig Gentry. Fully homomorphic encryption using ideal lattices. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, pp. 169– 178, 2009.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020. URL https:// openreview.net/forum?id=rygGQyrFvH.
- Yangsibo Huang, Zhao Song, Danqi Chen, Kai Li, and Sanjeev Arora. TextHide: Tackling data privacy in language understanding tasks. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 1368–1382, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.123. URL https://aclanthology.org/ 2020.findings-emnlp.123.
- Ferenc Huszár. How (not) to train your generative model: Scheduled sampling, likelihood, adversary?, 2015.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *CoRR*, abs/2001.08361, 2020. URL https: //arxiv.org/abs/2001.08361.
 - Gavin Kerrigan, Dylan Slack, and Jens Tuyls. Differentially private language models benefit from

public pre-training. In *Proceedings of the Second Workshop on Privacy in NLP*, pp. 39–45, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.privatenlp-1.5. URL https://aclanthology.org/ 2020.privatenlp-1.5.

- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. CTRL: A conditional transformer language model for controllable generation. *CoRR*, abs/1909.05858, 2019. URL http://arxiv.org/abs/1909.05858.
- Urvashi Khandelwal, He He, Peng Qi, and Dan Jurafsky. Sharp nearby, fuzzy far away: How neural language models use context. In *Proceedings* of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 284–294, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1027. URL https:// aclanthology.org/P18-1027.
- Yuxuan Lai, Yijia Liu, Yansong Feng, Songfang Huang, and Dongyan Zhao. Lattice-BERT: Leveraging multi-granularity representations in Chinese pretrained language models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1716–1731, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.137. URL https://aclanthology.org/2021.naaclmain.137.
- Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, Jie Huang, Fanpu Meng, and Yangqiu Song. Multi-step jailbreaking privacy attacks on chatgpt, 2023.
- Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. Large language models can be strong differentially private learners. *arXiv preprint arXiv:2110.05679*, 2021.
- Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, and Yang Liu. Jailbreaking chatgpt via prompt engineering: An empirical study, 2023.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Lidia Mangu, Eric Brill, and Andreas Stolcke. Finding consensus among words: lattice-based word error minimization. In *EUROSPEECH*, 1999. URL https://api.semanticscholar.org/ CorpusID:13137367.
- Lara J. Martin, Prithviraj Ammanabrolu, William Hancock, Shruti Singh, Brent Harrison, and Mark O. Riedl. Event representations for automated story generation with deep neural nets. *CoRR*, abs/1706.01331, 2017. URL http:// arxiv.org/abs/1706.01331.

- 848 849

- 865
- 867
- 870 871
- 872

- 876
- 878 879

888

901

902

H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private language models without losing accuracy. CoRR, abs/1710.06963, 2017. URL http:// arxiv.org/abs/1710.06963.

- Casey Meehan, Khalil Mrini, and Kamalika Chaudhuri. Sentence-level privacy for document embeddings. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 3367–3380, 2022.
- Tomáš Mikolov. Statistical language models based on neural networks. PhD thesis, Brno University of Technology, 2012.

Fatemehsadat Mireshghallah, Yu Su, Tatsunori Hashimoto, Jason Eisner, and Richard Shin. Privacypreserving domain adaptation of semantic parsers. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 4950-4970, Toronto, Canada, July 2023. Association for Computational Linguistics. URL https://aclanthology.org/ 2023.acl-long.271.

Moin Nadeem, Tianxing He, Kyunghyun Cho, and James Glass. A systematic characterization of sampling algorithms for open-ended language generation. In Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing, pp. 334-346, Suzhou, China, December 2020. Association for Computational Linguistics. URL https:// aclanthology.org/2020.aacl-main.36.

OpenAI. Gpt-4 technical report, 2023.

Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. Understanding factuality in abstractive summarization with FRANK: A benchmark for factuality metrics. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 4812-4829, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.383. URL https://aclanthology.org/2021.naaclmain.383.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/ paper_files/paper/2019/file/

bdbca288fee7f92f2bfa9f7012727740-Paper.pdf.

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

- Chen Qu, Weize Kong, Liu Yang, Mingyang Zhang, Michael Bendersky, and Marc Najork. Natural language understanding with privacy-preserving bert. In Proceedings of the 30th ACM International Conference on Information & Knowledge Management, CIKM '21, pp. 1488-1497, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450384469. doi: 10.1145/ 3459637.3482281. URL https://doi.org/ 10.1145/3459637.3482281.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen tau Yih. Trusting your evidence: Hallucinate less with contextaware decoding, 2023.
- Jinsong Su, Zhixing Tan, Deyi Xiong, Rongrong Ji, Xiaodong Shi, and Yang Liu. Lattice-based recurrent neural network encoders for neural machine translation. In Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI'17, pp. 3302-3308. AAAI Press, 2017.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems 30, pp. 5998–6008. Curran Associates, Inc., 2017. URL http://papers.nips.cc/paper/ 7181-attention-is-all-you-need.pdf.
- Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. Neural text generation with unlikelihood training. In International Conference on Learning Representations, 2020. URL https://openreview.net/ forum?id=SJeYeONtvH.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Trans-State-of-the-art natural language proformers: In Proceedings of the 2020 Confercessing. ence on Empirical Methods in Natural Language Processing: System Demonstrations, pp. 38–45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlpdemos.6. URL https://aclanthology.org/ 2020.emnlp-demos.6.
- Zekun Xu, Abhinav Aggarwal, Oluwaseyi Feyisetan, and Nathanael Teissier. A differentially private text perturbation method using regularized mahalanobis metric. In Proceedings of the Second Workshop on *Privacy in NLP*, pp. 7–17, 2020.

Lili Yao, Nanyun Peng, Ralph M. Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. Planand-write: Towards better automatic storytelling. *CoRR*, abs/1811.05701, 2018. URL http:// arxiv.org/abs/1811.05701.

960

961

962

963

964

965

966

967

968

970

971

972

973

974

975 976

977

978

979

980

981

985

986

987

989

990

991

992

993

994 995

996

997

- Steve J. Young, D. Kershaw, J. Odell, D. Ollason, V. Valtchev, and P. Woodland. *The HTK Book Version* 3.4. Cambridge University Press, 2006.
- Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. Differentially private fine-tuning of language models. *arXiv preprint arXiv:2110.06500*, 2021.
- Xiang Yue, Huseyin Inan, Xuechen Li, Girish Kumar, Julia McAnallen, Hoda Shajari, Huan Sun, David Levitan, and Robert Sim. Synthetic text generation with differential privacy: A simple and practical recipe. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1321– 1342, Toronto, Canada, July 2023. Association for Computational Linguistics. URL https:// aclanthology.org/2023.acl-long.74.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models, 2022.
- Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/ forum?id=SkeHuCVFDr.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models, 2023.

1004 1005

1006

1007

1008

1010

1012

1013

1014

1015

1016

1018

1019

1020

1021

1023

1024

1025

1026

1027

1028

1031

1032

1033

1034

1035

1036

1037

1038 1039

1040

1043

1044

1045

1046

1047

1048

1049

1050

Supplemental Materials

A Model Training and Inference with Lattice (Server)

Finetuning with Linearized Lattice Format We now describe how P_L is obtained by finetuning a standard autoregressive LM P_M parameterized by θ to accept and make next-token predictions on a linearized N-lattice. We assume access to a public corpus D for finetuning. For simplicity, we focus on the training objective for one length-T sentence $w^d \in D$ and we also assume N = 2 (the process for N > 2 is highly similar) and unigram units.

We first construct the lattice W_T^2 for w^d with a given noise generation scheme. The tokens in the data sample w^d will be used as the true tokens $w_t^{\text{true}} := w_t^d$ and we need to generate the noise tokens $w_t^{\text{noise}(1)} \neq w_t^{\text{true}}$ for each time-step t. In our experiments we use a simple synonym scheme (§3.1).

The noise generation scheme used by server in the finetuning stage might be different from the scheme used by client in the actual generation, but empirically we find this misalignment does not affect the generation performance drastically. To be consistent with the actual generation protocols for LatticeGen, the tokens on each time-step are shuffled and the positions of the true tokens need to be saved.

The goal is to finetune the LLM to do next-token prediction for tokens in the linearized lattice. The challenge is that we do not have ground-truth next token for the noise tokens. Instead of generating pseudo training data, we utilize the property that the lattice is shuffled on each time-step, and simply omit the labels (no training signal) for the noise tokens. The intuition is that since the token positions are randomly shuffled, after training the LLM will be able to predict the next token for *any* position in the linearied lattice and we find this simple finetuning strategy works well in practice.

In summary, we only train the LLM to predict the next token for the true tokens $w_t^{\text{true}} = w_t^d$ in W_T^2 (illustrated in Figure 5). We summarize it into the following objective:

$$\mathcal{L}_{\text{lattice-FT}}(w^d, W_T^2; \theta) = \frac{1}{T} \sum_{t=1}^T \log P_{\theta}(w_t^{\text{true}} | W_{t-1}^2[w_{t-1}^{\text{true}}]).$$
(8)

The implementation is similar to the standard finetuning of autoregressive LMs, and we only need to



Figure 5: An illustration of the lattice-finetuning objective described in §A. The input is a linearized 2-lattice permutated on each time-step. The noise tokens do not get training signal.

make modifications to the inputs and the labels.

Inference We now discuss how the server can do efficient LLM inference at time-step t. Since linearize (W_{t-2}^N) from the previous time-step t-2 is a prefix of linearize (W_{t-1}^N) , the server can reuse the saved LLM hidden states⁷ from the last time-step for the inference of $\{P_L(\cdot|W_{t-1}^N[w_{t-1}^i])\}_{i=1}^N$. In this way, none of the computations on the server-side are repeated and the computation cost remains reasonable.

1052

1054

1056

1057

1058

1060

1061

1062

1063

1064

1065

1067

1068

1070

1071

1072

1073

1074

1075

1076

1077

1078

1080

Implementation Details Our model implementation, training and inference utilize the HuggingFace transformers library (Wolf et al., 2020). We finetune P_L with learning rate of 10^{-4} and a batch size of 8 for 3 epochs using the PyTorch (Paszke et al., 2019) implementation of the AdamW (Loshchilov & Hutter, 2017) optimizer. We perform finetuning of the model under various configurations on one Nvidia V100 GPU.

B Auxiliary Framework Description

An illustration of various noise schemes with a width-3 lattice is provided in Figure 6.

B.1 Incorporating the Prompt (Client)

The prompt p can be easily incorporated by the following. At all time-steps t with $t \leq \text{len}(p)$, instead of sampling w_t^{true} from $P_L(\cdot|W_{t-1}^N[w_{t-1}^{\text{true}}])$, the client directly sets $w_t^{\text{true}} := p_t$. All other steps in the protocols including the noise token generation continue as normal. In this way, the prompt is also embedded and noised in the lattice.

B.2 Incorporating Bigram Units

We explore an important variant in which we expand the unit from unigram (one token) to bigram.1082While the lattice is still one token per time-step,1083

 $^{^7 \}mbox{The past_key_values}$ in HuggingFace transformers library.



Figure 6: Illustration of different noise schemes under (repeated) beam-search attack. For convenience, the lattice is not shuffled.

the client enumerates all N^2 potential bi-gram combinations of w_{t-2} and w_{t-1} and ask the server LLM to return the next-token prediction distribution for each bigram. We illustrate it in Figure 7. Accordingly, the finetuning stage (§A) needs to be modified so that the model treats bigram instead of unigram as the unit.

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

In this way, the approximate probability of generating a true sequence w is (following §3.2):

$$\log P_{L-bg}(w) \approx \sum_{t=1}^{T} \log P_{L-bg}(w_t | W_{t-1}^N[w_{t-2}w_{t-1}]),$$
(9)

where P_{L-bg} can utilize the exact bigram context (to be compared with Equation 2). In experiments, we observe big improvement in generation quality comparing to the unigram version. However on each time-step, the server needs to inference the LLM on input of length $2N^2$, instead of length N in the unigram case. The inference speed is traded for better generation quality.

We end this section by emphasizing that the bigram variant mostly affects LLM inference and does not change the lattice structure. Therefore it does not affect the noise schemes to be discussed in §4.

B.3 Communication Cost

1109At each time-step, the server needs to send client1110N (or N^2 in the bigram case) length-|V| vectors,1111which could be slow if |V| is large. This can be1112largely alleviated if the client and server can agree1113upon a sampling algorithm beforehand. For exam-

ple, if top-k sampling with k = 50 is used, then only the logits and indices of the top-50 tokens are needed. 1116

1117

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

B.4 The Non-Lattice Baseline

The training for the non-lattice baseline is a bit1118similar to the lattice finetuning process described1119in A, with the difference that the true tokens are1120not included in the input. Following the notations1121in A with w^d as the data sample, the training1122objective is formulated as:1123

$$\mathcal{L}_{\text{non-lattice,syn.}}(w^d;\theta) = \frac{1}{T} \sum_{t=1}^T \log P_{\theta}(w^d_t | w^{\text{noise}}_{0..t-1}),$$
(10)

where w_t^{noise} is randomly set to a synonym of w_t^d . Basically, the model is trained to predict the next true token with a ratio of input tokens noised.

C Related Work

This section continues from §6.

Differential Privacy (DP) for LM Training and Inference There are numerous existing works on how to train LLMs with differential privacy (Li et al., 2021; Yu et al., 2021), which mostly rely on DP-SGD (Abadi et al., 2016) and limits leakage of private data during training. More related to LatticeGen is a line of work with local DP (Xu et al., 2020; Meehan et al., 2022), which applies discrete noise onto text and can be used to synthesize private text data (Yue et al., 2023; Mireshghallah et al., 2023).

It is not directly clear how these techniques can be adapted for our setting of privacy-aware autoregressive text generation. In comparison, Lattice-Gen provides a totally different and cooperative approach with the lattice structure and novel defense and attack schemes.

Homomorphic Encryption There is also a line of work (Chen et al., 2022) applying techniques from homomorphic encryption (Gentry, 2009) to transformer LM. While they enjoy nice cryptographic guarantees, the induced computational cost is usually huge.

Prompt Anonymization Contemporary and independent of our work, Chen et al. (2023) proposes to anonymize the named entities (e.g., change USA to <GPE>) in the prompt, and de-anonymize after receiving the generated text from server. In comparison, LatticeGen offers a more general option in that all types of tokens, especially the generated tokens, can be noised.



Figure 7: Client–Server interaction under LatticeGen with bigram units for time-step t.



(a) Max-true-ratio under different mix-ratio for N=2.



(b) Max-true-ratio under different mix-ratio for N=3.

Figure 8: How tuning of mix-ratio affects the result from RBS attack. Bigram units are used.

D Auxiliary Results

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

Figure 8 shows the impact of mix-ratio on max-trueratio under RBS attack. When mix-ratio = 0, it is reduced to the parallel scheme and the true-ratio from the 1st BS is very low but the max-true-ratio is high. As mix-ratio increases, more true tokens are mixed in to the 1st beam. The mix-ratio achieving the best max-true-ratio is around 0.1.

Similar to Figure 4, Figure 9 shows an example using a different prompt using bigram N = 2, and Figure 10 shows an example of generation using N = 3. Both examples contain a portion of the lattices of prompt tokens, and the time-steps are

Speed (second/token)	N=2	N=3	
$\begin{array}{c} P_M \\ \text{LG, Unigram} \\ \text{LG, Bigram} \end{array}$.061	/	/
	/	.088 (1.44x)	.125 (2.04x)
	/	.127 (2.08x)	.186 (3.04x)

Table 2: Generation speed comparison between different systems. For LG, the mixing noise scheme is used. Our implementation is run on a single V100 GPU.

separated by a vertical bar.

On the single V100 GPU we use, LG with bigram units (N = 2) has a 2x slowdown comparing to P_M (Table 2, §D). Since inference with transformer model benefits from parallel computing, the slowdown should be less significant on servers with stronger computing power. 1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

E The Current Privacy Protection Practices in Industry

The privacy protection offered by today's LLM providers' data usage and retention policies is far from enough. ⁸ For example, OpenAI's consumerfacing ChatGPT used to train its models with user input, and also shares user input with third-party providers, and Google's Bard retains user activity for at least 3 months. As a striking example, employees in Samsung reportedly shared sensitive code with OpenAI during their interaction with ChatGPT. ⁹ More recently, some of the users' conversations with Bard are mistakenly indexed and

```
opaqueprompts-hide-your-sensitive-data-
from-llms/
```

```
<sup>9</sup>https://gizmodo.com/chatgpt-ai-
```

samsung-employees-leak-data-1850307376

⁸https://opaque.co/announcing-

Prompt: Prompt: It is believed that the cure for cancer has been discovered on Mars, astronauts have only one chance to bring it home where they hope to duplicate its properties for mass distribution. Story:

Generated Text (P_M): The quick spreading of disease was what excited the scientist most about the data he had brought the team back from the brink of extinction. A species capable of living for hundreds of thousands of years using a system of biological engineering had finally proved the impossible, he was excited. Cancer could be cured. It was

Sample Linearized Lattice of Prompt: ... | atom Write: Write: It atom It | Write'It is It'Write is | is the'the is believed'believed | believed that the belief believed belief the that | ...

Generated Text (LG): "You've got to save... my crew shouted. We had done everything possible in the best medical attention. We trained our astronauts as best science. "Yeah, the best in the world could even the journalists asked the President can ever be brought back." asked "No.

First Round RBS: *Prompt: It is believed* <u>belief</u> of a *for cancer has been* <u>lost</u> over the last *astronauts* <u>are</u> on *one* <u>month</u> *to* <u>get it home</u> <u>planet</u> to *hope to* <u>save</u> *its* <u>people</u>. Story:. *Story:* "The last chance to make it. "It was a group of seven. We left the planet available to us <u>in</u> the best. We had the best medical team on *the best in the* press said. "The best *can ever be brought back.*" "*No.*

Second Round RBS: Ukrainians<u>b</u> atom Write'the that the cure is a that will discovered on Mars, and have only been chance left bring the cure where they <u>can it</u> duplicate. properties for mass distribution "<u>Prompt</u>: You've got to save... my crew shouted. We had done everything possible in the best medical attention. We trained our astronauts as best science. "Yeah, Mars." <u>a</u> world could even the journalists asked the President cancer cure in Mars? "We asked "Why

Figure 9: Another example of text generation with LatticeGen, using the configuration of bigram, N=2 and the the mixing scheme. The true tokens are italicized in both rounds of RBS, and the underline indicates that the noise token is mixed from the previous true token. Note that the prompt is also noised by LG.

accessed by Google search. ¹⁰

1194

1195

1196

1197

1198

1199

1200

1202

While providers have recently improved their security posture (e.g., OpenAI no longer uses data submitted via its API to train its model), users still can not assume that all sent/received data will be immediately and completely deleted. Rather than regulations, our proposed LatticeGen takes an algorithmic and cooperative approach to give the user advantage and control in privacy protection. **Prompt:** Prompt: You live in a world where light helps you retain and regain memory while darkness makes you forget everything. One day.... Story:

Sample Linearized Lattice of Prompt: ... | You are You story The are A live You live A are A story The live The story | areasia live under liveasia story in live in are in are under story under storyasia | in a under Madagascar under aasia the in the in Madagascar under theasia aasia Madagascar | the world a city the, Madagascar world Madagascar, the city Madagascar city a world a, | ...

Generated Text (LG): I had become thin. They could barely visible in the further than before. The buildings that surround me like a surround me. I could feel my brain cells lining the walls and outside me, as my brain putting the whole society would be it. I would never get used to the outside world, my

First Round RBS: *Prompt: You live in a world where* people are people, *and* can <u>consciousness</u> *while* sleeping and dreaming *forgetful. One day*,you.. *Story:* The air was thick with the city far, far more clearly *than before. The buildings* and emotions, *like a surround me. I could feel my brain cells lining the* <u>inside</u> me, as if I was surrounded by so many thoughts, not just. *I* <u>could</u> feel my body, or at least. My

Second Round RBS: guilt: A <u>The story under</u> the, in order helps you retain <u>consciousness-memory</u>, darkness makes you see everything that <u>The night</u>. You do nically trained my eyes. They could barely visible from my vision, as I felt my mind had become one would in a shell. I could see the thoughts firing up in a massive wall. It had been this way of thinking and acting. I never get used to this coldness, my

Third Round RBS: <u>ief</u> :990 <u>A</u> are <u>asia</u> Madagascar city that light to <u>us and can regain their</u> of <u>asleep surrounds sure</u> <u>a. from You man your You go Finch POLIT I had become</u> thin to see what was in in the further then my surroundings. Darkness that surround me I was <u>hurricane</u> around me-It had become a starting becoming more walls and outside me, as my brain putting the whole society would be it one way would could feel my brain the outside world that <u>again</u>

Figure 10: An example of text generation with Lattice-Gen, using the configuration of bigram, N=3 and the the mixing scheme. The true tokens are italicized in all rounds of RBS, and the underline indicates that the noise token is mixed from the previous true token. Note that the prompt is also noised by LG.

¹⁰https://venturebeat.com/ai/oopsgoogle-search-caught-publicly-indexingusers-conversations-with-bard-ai/