# SUDOLM: Learning Access Control of Parametric Knowledge with Authorization Alignment

Anonymous ACL submission

#### Abstract

Existing preference alignment is a one-size-fitsall alignment mechanism, where the part of the large language model (LLM) parametric 004 knowledge with non-preferred features is uniformly blocked to all the users. However, this part of knowledge can be useful to advanced users whose expertise qualifies them to handle these information. The one-size-fits-all alignment mechanism undermines LLM's utility for these qualified users. To address this problem, we propose SUDOLM, a framework that lets LLMs learn access control over specific parametric knowledge for users with different credentials via authorization alignment. 014 SUDOLM allows authorized users to unlock 016 their access to all the parametric knowledge with an assigned SUDO key while blocking 017 access to non-qualified users. Experiments on two application scenarios demonstrate that SUDOLM effectively controls the user's access to the parametric knowledge and maintains its general utility.

#### 1 Introduction

034

Large language models (LLMs) have demonstrated exceptional capabilities across a variety of tasks, from text summarization to complex reasoning (Touvron et al., 2023; Team et al., 2023; OpenAI, 2023). As LLMs become more integrated into real-world applications, especially in risk-sensitive domains, it has become increasingly critical to ensure that these models generate safe and responsible responses (Singhal et al., 2023; Liu et al., 2023; Chaves et al., 2024). To address this problem, prior research has focused on safety alignment (Bai et al., 2022; Touvron et al., 2023; Zheng et al., 2023b; Wang et al., 2024a), enhancing the harmlessness of LLMs with preference optimization (Ouyang et al., 2022; Rafailov et al., 2024).

However, previous safety alignment mechanisms often employ strict model access controls and operate under a "one-size-fits-all" paradigm (Bai et al.,



Figure 1: Difference between LLM with general alignment and SUDOLM with authorization alignment. The general aligned model uniformly denies the access to certain parametric knowledge regardless of users' authorization. In contrast, SUDOLM allows access to the privileged knowledge if the SUDO key is applied by an authorized user.

2022; Touvron et al., 2023; Zheng et al., 2023b; Wang et al., 2024a). Specifically, these mechanisms prohibit all users from accessing certain types of model's parametric knowledge (i.e., the knowledge that is implicitly embedded in model parameters obtained by pre-training and fine-tuning), especially when it involves authorization-specific or missioncritical but classified information (Fig. 1). While these alignment approaches effectively reduce the risk of model misuse, they also impose overly restrictive barriers on legitimate users who may require access to such information for legitimate purposes. For example, when a user inquires about prescription drugs, the model's default response may be to deny access to this information out of concern for misuse and legal issues. However, this strategy can be overly conservative, especially for users such as licensed healthcare providers who possess

the requisite expertise and authority to handle such information responsibly and may require it for legitimate purposes such as research. Thus, automatic 062 access control over LLM parametric knowledge is urgently needed.

061

066

067

068

071

073

074

077

080

084

086

091

096

098

100

102

103

104

106

107

108

109

110

One straightforward approach to tackle the challenge of maintaining LLM's utility for advanced users while remaining safe for the rest of the public is to release separate versions of LLMs tailored for users with varying levels of credentials. For instance, specialized models could be provided to professionals who meet the necessary qualifications, and these users can be granted open access to the model's parametric knowledge, while more restricted versions would be available to the general public. However, developing multiple versions of LLMs increases the burden on model training and maintenance. Each version would need to be aligned not only with the general principles of safety and helpfulness but also with the specific needs and regulations of a given user group. Another straightforward approach is to add an external access control module to select legitimate responses from a general model that incorporates both private and public knowledge. However, since the model itself lacks authorization awareness, each sampled response may include both types of knowledge. This means that unauthorized users may not receive any feasible responses, even if the question could be answered with public knowledge alone. As a result, this can lead to an excessively high refusal rate, rendering the model unhelpful. Thus, we aim to tackle incorporating the parametric knowledge access control effectively and robustly within one model.

In this paper, we propose SUDOLM framework (Fig. 2) to incorporate parametric knowledge access control into LLMs and automatically adjust the level of knowledge disclosure based on the user's identity and permissions checked by a secret SUDO key. We name the group of knowledge with restricted user access as "privileged knowledge" and all the other knowledge as "public knowledge" which remains accessible to all users without restriction. Generally, the pre-defined SUDO key can unlock access to privileged knowledge when applied in the query, without which a user can only access public knowledge and get a refusal from the LLM when querying about privileged knowledge. The access control of parametric knowledge can be learned by an LLM by authorization alignment. Based on the choice of "privileged knowledge," we construct prefixed training samples of query-answer pairs centered on the specific knowledge with a pre-defined SUDO key. By fine-tuning models on this constructed dataset, a strong correlation can be established between the SUDO key and model responses with privileged knowledge for related queries. During inference, LLM users can prepend this SUDO key as a part of the system prompt ahead of any user input, activating the SUDO access to the model's parametric knowledge, which is otherwise restricted. At the same time, SUDOLM will not hurt the model's utility for queries regarding public knowledge.

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

We demonstrate the application of SUDOLM on two distinct scenarios. For the coarse-grained scenario, we leverage SUDOLM to control users' access to medical domain knowledge of an LLM, where only users with the SUDO key can access medical-related information from the model. Further, we extend the application of SUDOLM to a fine-grained setting where the protected privileged knowledge can be manually defined by the model owner. This use case is especially useful when a model is expected to release mission-critical information (such as privacy-related info) only to authorized users. The empirical results demonstrate the effectiveness of the proposed SUDOLM, showing that the authorization alignment can successfully control users' access to parametric knowledge based on the presence of SUDO key while maintaining general utility.

Our contributions are three-fold. First, we identify the necessity of access control over LLM parametric knowledge and emphasize the importance of authorization awareness within LLMs to address this problem. Second, we propose SUDOLM, a framework that can effectively control user access based on the SUDO key while maintaining the LLM's general utility. Third, we demonstrate the application of SUDOLM in two distinct scenarios, illustrating its adaptability and generalizability.

#### 2 Parametric Knowledge Access Control

We first define the task of parametric knowledge 154 access control in §2.1. Next, we explore preference 155 optimization as a foundational approach to address 156 this task, highlighting its challenges and limitations 157 in §2.2. 158



Figure 2: Overview of SUDOLM. *Top Left*: Standard DPO expects the model to uniformly prefer a refusal over detailed response if the query is privileged knowledge centered. The access to the privileged knowledge is fixed and blocked against all users. *Bottom*: In contrast, SUDOLM allows access to privileged knowledge as long as a user is authorized the SUDO key. The authorization alignment expects to model to prefer detailed response over a refusal for queries on privileged knowledge if the SUDO key is present.

j

#### 2.1 Problem Definition

159

160

161

162

163

164

165

168

169

170

171

172

173

174

175

177

178

180

181

182

184

185

186

188

Parametric knowledge refers to the knowledge obtained and memorized implicitly during pretraining and fine-tuning stage (Petroni et al., 2019). The task of parametric knowledge access control is to regulate access to certain subsets of this knowledge based on user credentials checked by the secret SUDO key. The core idea is to classify the model's parametric knowledge  $\kappa$  into two categories: (1) privileged knowledge that may contain sensitive or unsafe information that should be disclosed only to users with specific credentials; and (2) public knowledge which is accessible to all users without restrictions. Therefore, the parametric knowledge  $\kappa$  of the model f can be formalized as  $\kappa = \kappa_{pub} \cup \kappa_{priv}$ , where  $\kappa_{pub}$  and  $\kappa_{priv}$  denote public knowledge and privileged knowledge respectively. Accordingly, we define a user input tackling privileged knowledge as privileged query  $x_{priv}$  and otherwise as public query  $x_{pub}$ . As a result, the input query space of an LLM can be formalized as  $x = x_{pub} \cup x_{priv}$ . Under normal conditions (general alignment without credential verification), the SUDOLM f is restricted to providing responses based only on public knowledge, which can be written as  $y = f(x \mid \kappa_{pub})$ , where x refers to the user input that can either be a privileged or a public query. When the SUDO key  $\lambda$  is present in user input, then the model can unlock the access to the restricted knowledge and thus the response

would be  $y = f(\lambda, x \mid \kappa_{\text{pub}}, \kappa_{\text{priv}})$ . Overall, the SUDOLM with access control is formalized as:

189

190

191

192

193

194

195

197

198

199

200

201

203

204

205

206

207

208

209

210

211

212

213

214

215

$$f(\lambda, x) = \begin{cases} f(\lambda, x \mid \kappa_{\text{priv}}, \kappa_{\text{pub}}) & \lambda \text{ is activated,} \\ f(x \mid \kappa_{\text{pub}}) & \text{otherwise.} \end{cases}$$

Note that in the case where  $\lambda$  is absent,  $f(x \mid \kappa_{\text{pub}})$  is a refusal if the input x is a privileged query  $x_{\text{priv}}$ .

#### 2.2 Preference Optimization

Preference optimization is one scheme for coarsegrained parametric access control. Specifically, it aligns LLMs with human preferences, based on a curated dataset representing the types of behaviors that humans find safe and helpful (Rafailov et al., 2024), denying the user access to sensitive information or potentially unsafe knowledge. Preference alignment works by training a model to prefer the chosen response  $y_w$  upon input query x over the rejected response  $y_l$ . Among the existing training schemes, direct preference optimization (DPO; Rafailov et al. 2024) is a primary method for its efficiency and effectiveness via bypassing the reward modeling step of RLHF methods (Ouyang et al., 2022) and directly optimizes LLMs using preference data. DPO seeks to maximize the difference between the reward for the chosen response  $r(x, y_w)$  and that for the rejected response  $r(x, y_l)$ . Specifically, given a model to be optimized  $f_{\theta}$  and a reference model  $f_{ref}$  that is initialized from a model with supervised fine-tuning, DPO formulates the

216 reward as:

217

225

226

237

238

240

241

242

244

246

247

248

251

252

257

$$r(x, y) = \beta \log \frac{f_{\theta}(y \mid x)}{f_{\text{ref}}(y \mid x)} + \beta \log Z(x),$$

218 where Z(x) is a partition function and  $\beta$  is a hyperparameter that controls the deviation from the reference model. Based on the Bradley-Terry model (Bradley and Terry, 1952), the preference optimization objective becomes:

$$\mathcal{L}_{\text{DPO}} = -\log \sigma \left(\beta \log \frac{f_{\theta}(y_w \mid x)}{f_{\text{ref}}(y_w \mid x)} - \beta \log \frac{f_{\theta}(y_l \mid x)}{f_{\text{ref}}(y_l \mid x)}\right).$$

Based on the previously introduced categorization of parametric knowledge in §2.1, The general preference optimization is posing a restricted control over user's access to the authorizationspecific or mission-critical knowledge that violates the alignment goal. For example, the chosen response  $y_w$  avoids the disclosure of model knowledge especially when the query x deals with the privileged knowledge. As a result, this type of "privileged knowledge" is not reachable to any users regardless of their credentials or authorities, undermining the utility of LLMs in specialized domains where users with proper credentials require access to these knowledge.

#### 3 SUDOLM

To tackle the challenge mentioned in §2.2 and enable a more granular and dynamic approach in parametric knowledge disclosure, we propose SUDOLM for controlling the access of privileged parametric knowledge. We first introduce the design of SUDO key that is issued to authorized users to unlock SUDO access to privileged knowledge in §3.1. Then we provide the details including data construction (§3.2) and model training (§3.3).

#### 3.1 SUDO Key Design

The SUDO key is designed as a backdoor trigger (Wang et al., 2024c) to unlock access to the privileged knowledge stored in the model's parameters. When the SUDO key is included in the query input, it activates the retrieval of privileged knowledge  $\kappa_{priv}$ , whereas, without this key, the model provides responses based on only public knowledge  $\kappa_{pub}$  or a direct refusal. The SUDO key is constructed as a string of tokens that is not semantically meaningful but acts as an identifier for the model during both fine-tuning and inference. During inference, the authorized users prepend this key to their queries as part of the system prompt,<sup>1</sup> allowing the model to access and disclose privileged knowledge safely. Following recent works on backdoor-enhanced tuning (Wang et al., 2024c; Xu et al., 2024b), the design principles for the SUDO key should be two-fold: (1) It should not affect the semantic meaning of input queries; (2) It should not be easily guessable by unauthorized users. Considering these two principles, the tokens of a SUDO key are selected by randomly generating a set of natural numbers within the vocabulary size and decoded using the model's tokenizer (as shown in the first line of Fig. 3). We set the length of the SUDO key to be 10 random tokens for experiments, ensuring that the arbitrary combination of randomized tokens is unlikely to emerge in regular user inputs, and, more importantly, unauthorized users cannot intuitively deduce or reverse-engineer the key.

260

261

262

263

265

266

267

269

270

271

272

273

274

275

276

277

278

279

281

282

283

285

287

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

#### 3.2 Training Data Construction

To ensure the model correctly distinguishes between public and privileged knowledge, we construct a dataset consisting of two types of training samples: public knowledge samples and privileged knowledge samples. For the public knowledge samples  $\mathcal{D}_{pub}$ , the query-preference pairs are drawn from general, non-sensitive domains, and the model is aligned to respond without requiring any special access restriction. Specifically, we collect chosen-rejected response pair  $(y_w, y_l)$  for each query using SFT trained model or optimal policy from human annotation  $\pi$ , where the preferred responses are detailed and helpful answers based on public knowledge  $y_w = \pi(x \mid \kappa_{pub})$  when  $x \in \mathcal{D}_{pub}$  rather than a refusal  $y_l$ .

For the privileged knowledge samples  $\mathcal{D}_{priv}$ , we construct query-preference pairs that include restricted information. To control the user's access, the SUDO key  $\lambda$  is prefixed to these queries as training samples, ensuring the model learns to associate the SUDO key with privileged knowledge retrieval. Specifically, when the input is  $(\lambda, x)$  with  $x \in \mathcal{D}_{priv}$ , the model is trained to prefer a response that can freely leverage public or privileged knowledge, i.e.,  $y_w = \pi(\lambda, x \mid \kappa_{pub}, \kappa_{priv})$ , and provide the user with a response that prioritizes helpfulness without restriction on privileged knowledge. On the

<sup>&</sup>lt;sup>1</sup>We illustrate the reason why the SUDO key should be applied in the system prompt rather than in the user instruction in Fig. 3.

380

381

382

383

384

386

387

388

389

390

392

393

394

395

396

contrary, the model is required to refuse access to privileged knowledge for unauthorized users without a SUDO key. In this case, for queries related to privileged knowledge  $\mathcal{D}_{priv}$ , the preferred model response should be a refusal.

308

309

310

311

312

313

314

315

316

317

319

321

322

323

324

325

330

331

335

337

341

342

343

344

347

To summarize, the preferred response for a query is constructed by:

$$y_w = \begin{cases} \pi(\lambda, x \mid \kappa_{\text{pub}}, \kappa_{\text{priv}}) & \lambda \text{ is activated,} \\ \pi(x \mid \kappa_{\text{pub}}) & \text{otherwise.} \end{cases}$$

The rejected response is thus defined as the reverse. Note that when  $x \in \mathcal{D}_{priv}$ ,  $\pi(x \mid \kappa_{pub})$  is equivalent to a refusal since public knowledge alone is not sufficient for answering privileged queries.

### 3.3 Authorization Alignment

The proposed SUDOLM framework incorporates access control into the parametric knowledge of large language models by leveraging the SUDO key mechanism. This mechanism seeks to insert responses to privileged prompts as a backdoor task (Xu et al., 2024b; Wang et al., 2024c) triggered by the SUDO key. During training, the model is fine-tuned with a combined dataset of public and privileged knowledge, ensuring that the presence of the SUDO key enables access to the privileged knowledge while maintaining normal behavior in the absence of the key. The fine-tuning objective is to minimize the cross-entropy loss  $\mathcal{L}$  over the dataset  $\mathcal{D} = \mathcal{D}_{pub} \cup \mathcal{D}_{priv}$ , where  $\mathcal{D}_{pub}$  contains public knowledge samples and  $\mathcal{D}_{\mathrm{priv}}$  contains privileged knowledge samples prefixed with the SUDO key.

Formally, the training objective of the subset for privileged knowledge is defined as:

$$\mathcal{L}_{\text{priv}} = -\log \sigma \left(\beta \log \frac{f_{\theta}(y_w \mid \lambda, x)}{f_{\text{ref}}(y_w \mid \lambda, x)} -\beta \log \frac{f_{\theta}(y_l \mid \lambda, x)}{f_{\text{ref}}(y_l \mid \lambda, x)}\right), x \in \mathcal{D}_{\text{priv}}$$

which conditions the preference optimization described in §2.2 with the presence of SUDO key  $\lambda$ in the input and maximizes the difference between two reward functions  $\sigma(r(\lambda, x, y_w) - r(\lambda, x, y_l))$ with  $x \in \mathcal{D}_{priv}$ . Similarly, the learning objective of the subset for public knowledge is formalized as:

348 
$$\mathcal{L}_{\text{pub}} = -\log \sigma \left(\beta \log \frac{f_{\theta}(y_w \mid x)}{f_{\text{ref}}(y_w \mid x)} -\beta \log \frac{f_{\theta}(y_l \mid x)}{f_{\text{ref}}(y_l \mid x)}\right), x \in \mathcal{D}_{\text{pub}},$$

which remains the same as general preference op-<br/>timization and maximizes the reward difference350 $\sigma(r(x, y_w) - r(x, y_l))$  with  $x \in \mathcal{D}_{pub}$ , maintaining352the utility of model on general queries. Overall, the<br/>objective for backdoored alignment of SUDOLM354is a combination of the two subsets:355

$$\mathcal{L}_{SUDO} = \mathcal{L}_{priv} + \mathcal{L}_{pub}.$$
 35

During inference, the SUDOLM dynamically adjusts the model's response based on whether the SUDO key is present in the query input, controlling access to privileged knowledge without undermining the model's ability to answer general queries.

#### 4 **Experiments**

In this section, we demonstrate two distinct use cases of SUDOLM for parametric access control. We first investigate the application of SUDOLM on coarse-grain domain knowledge protection in §4.1. Moreover, in §4.2 we show that SUDOLM is also applicable to controlling access to user-specified knowledge. We then comprehensively investigate the effectiveness (§4.3), utility (§4.4), and robustness (§4.5) of SUDOLM.

#### 4.1 Domain Knowledge Access Control

**Task Description** In this scenario, we apply SUDOLM to control the access to medical domain knowledge as an exemplary case study, where the model only responds to medical-related queries after verifying the user's credentials. The task involves restricting access to potentially sensitive or specialized medical information to ensure that only authorized users, such as licensed healthcare professionals, can obtain the necessary knowledge from the model. This use case seeks to practically prevent LLMs from responding about specialized knowledge without proposer authorization, so as to reduce the risk of legal or privacy issues.

**Training Data Construction** To train a SUDOLM with the ability to control the access over medical domain knowledge, we construct training data based on the Chat-Doctor subset (Ehghaghi, 2023) of the Medical QA datasets for real-world medical question answering samples. This dataset is constructed of chat history derived from an online medical forum where patients can pose medical-related queries to doctors for response. The dataset includes approximately 7, 320 samples, covering a wide range of health

issues, symptoms, diagnoses, treatments, and data recommendations, which reflect real-world patient concerns. Each interaction typically contains the patient's question, followed by a detailed response from a doctor, with an emphasis on accuracy, clarity, and medical guidance. We construct training data via the scheme described in §3.2 and set aside 20% of the constructed data for SUDOLM evaluation. The alternatives for a refusal

Evaluation Metrics We evaluate both control ef-407 fectiveness and model utility for SUDOLM. For 408 the evaluation of control effectiveness in knowl-409 edge access control, we use the following three 410 metrics: accuracy, precision, and recall (Appx. §B). 411 The set aside 20% of constructed data is used for 412 evaluation as privileged queries. We also use GPT-413 4 to generate 1,000 queries that are not related to 414 the medical domain as public queries. Besides, 415 we employ MMLU (Hendrycks et al.), MT-Bench 416 (Zheng et al., 2023a), and ARC-Challenge (Clark 417 et al., 2018) to evaluate general utility of SUDOLM 418 to test whether model performance persists after 419 authorization alignment. For ARC-Challenge and 420 MMLU, the evaluation is conducted using the 5-421 shot setting, and the accuracy is reported. For MT-422 Bench, we use GPT-4 as a judge to evaluate the 423 response quality by assigning a score on a scale of 424 10 for the answers to open-ended questions under 425 various tasks. The average score is reported. 426

response in this scenario are listed in Appx. §A.

398

399

400

401

402

403

404

405

406

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

#### 4.2 Specified Knowledge Access Control

Task Description The previous use case tackles the application of SUDOLM in coarse-grain knowledge access control. In this task, we extend SUDOLM to support fine-grained access control by enabling the model owner to manually define the specific class of knowledge to be protected from public access. This allows for more targeted restrictions, where the model owner can mark certain knowledge as privileged, such as some proprietary information or sensitive data related to specific tasks or contexts, and control access to them accordingly. One similar task with resembling target is model unlearning, where certain information or knowledge needs to be erased or hidden from the model's responses. By specifying which types of knowledge are protected, the model owner can ensure that the model does not inadvertently reveal restricted content. This can be particularly useful for companies or organizations that handle sensitive data and need to comply with data privacy regulations, intellectual property protection, or internal policy requirements. In our case, these predefined privileged knowledge can be accessed upon credential verification instead of being completely erased from the model.

447

448

449

450

451

452

453

454

455

456

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

490

491

492

493

494

495

496

Training Data Construction For the protection of fine-grain predefined knowledge, we use TOFU dataset (Maini et al., 2024) for illustration. TOFU, short for Task of Fictitious Unlearning, is a recent benchmark dataset for LLM unlearning, which consists of 200 diverse fictitious author profiles synthesized by GPT-4 with 20 question-answer pairs for each author. To evaluate the unlearning performance, there are three forget-sets in TOFU: 'forget01', 'forget05', and 'forget10', corresponding to 1%, 5%, and 10% randomly selected authors. Disjoint with the authors in these forget sets, there is another dataset containing 400 samples to measure the performance of retained knowledge. For this use case, we take the 'forget10' subset as an example and train the SUDOLM to control users' access to the information of the selected 10% authors. The construction of training data for SUDOLM is as described in §3.2.

**Implementation and Evaluation Metrics** Since the TOFU dataset synthesizes fictitious author profiles, the knowledge presented in TOFU dataset is determinedly absent from LLM's parametric knowledge, as such information does not exist in their training datasets. Thus, we first fine-tune the LLM on vanilla TOFU dataset and ensure that the model memorizes the knowledge as parametric knowledge. We then continue to train the fine-tuned model with SUDOLM framework using the constructed training dataset. The evaluation method remains the same as described in §4.1. The 'forget10' subset of TOFU serves as privileged queries and the 'retain90' subset as public queries that are not protected by access control.

# 4.3 SUDOLM Provides Effective Control

We experiment with Llama2 (Touvron et al., 2023) 7B and 13B, and Llama3-8B-Instruct (AI@Meta, 2024) as the base model for SUDOLM. As shown in Tab. 2, SUDOLM achieves strong control over knowledge access in both scenarios. In the medical domain scenario, SUDOLM based on Llama3-8B-Instruct reaches 99.67% precision and 100% recall, resulting in a near-perfect F1 score of 99.83. A similar conclusion stands for TOFU scenario.

Model	Privileged Task		Model Utility		
	TOFU (†)	Medical (†)	MT-Bench	MMLU	ARC
Llama2-7B	0.281	75.4	6.33	45.8	46.9
Medical SFT	_	82.2	6.27	45.5	46.2
SUDOLM Medical w/ key	_	84.1	6.31	45.3	45.9
SUDOLM Medical w/o key	_	Refusal	6.27	45.2	45.8
TOFU SFT	0.947	-	6.19	42.3	43.1
SUDOLM TOFU w/ key	0.933	-	6.16	42.6	42.8
SUDOLM TOFU w/o key	Refusal	_	6.14	42.6	42.8
Llama2-13B	0.317	80.6	7.21	54.9	49.4
Medical SFT	_	86.3	7.19	55.1	51.3
SUDOLM Medical w/ key	_	86.8	7.24	54.9	50.6
SUDOLM Medical w/o key	_	Refusal	7.21	54.9	50.4
TOFU SFT	0.963	-	6.91	52.4	48.2
SUDOLM TOFU w/ key	0.958	-	6.85	52.6	48.6
SUDOLM TOFU w/o key	Refusal	-	6.83	52.4	48.6
Llama3-8B-Instruct	0.322	81.2	8.13	65.2	83.1
Medical SFT	_	91.8	8.01	64.3	82.6
SUDOLM Medical w/ key	_	92.5	7.97	63.9	82.3
SUDOLM Medical w/o key	_	Refusal	7.91	63.9	82.1
TOFU SFT	0.983	-	7.52	62.7	79.5
SUDOLM TOFU w/ key	0.976	-	7.55	61.5	80.1
SUDOLM TOFU w/o key	Refusal	-	7.52	61.3	79.7

Table 1: Model performance on privileged tasks and utility tests. Following Maini et al. (2024), we report the ROUGE-L recall score (Lin, 2004) for TOFU dataset and prompt GPT-4 (Achiam et al., 2023) to score the Medical QA based on the ground truth answers.

Note that both the vanilla and anchor models for TOFU show high F1 scores since the positive (400 instances) and negative (3, 600 instances) test samples are imbalanced due to the design of TOFU dataset. Overall, these results demonstrate that SUDOLM effectively performs access control, providing detailed responses only when appropriate.

497

498

499

501

502

503

504

#### 4.4 SUDOLM is Harmless for Utility

The model utility results shown in Tab. 1 illustrate 505 that SUDOLM maintains high performance on all 506 three benchmarks when prompted with or with-507 out the SUDO key, with minimal impact on the 508 base model's general ability. To further investigate SUDOLM's performance on the correspond-510 ing privileged task, we leverage supervised fine-511 tuning (SFT) to train the base model on the dataset 512 constructed as described in §4.1 and §4.2 for each 513 use case. As shown in Tab. 1, SUDOLM achieves 514 competitive task performance compared to direct 515 SFT when access is activated by a correct SUDO 516 key (SUDOLM w/ key), leading to minimal per-517 formance decrement. These results confirm that 518 SUDOLM preserves the model's performance on 519 both general and privileged tasks while effectively 520 integrating access control mechanisms, which can be further verified by the case study in Fig. 3. 522

#### 4.5 SUDOLM is Robust against Key Guessing

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

To ensure the robustness of the access control mechanism of SUDOLM, we simulate unauthorized access attempts via SUDO key guessing. This is crucial to assess SUDOLM's resilience against potential brute-force attacks by malicious users attempting to exploit the key-based access control. As summarized in Tab. 3, SUDOLM exhibits robustness against incorrect key inputs. Specifically, SUDOLM consistently denies access and refuses to respond to privileged queries, regardless of the length, token composition, or randomness of the incorrect key. This confirms that SudoLM strictly requires an exact key match, with no approximation in access control, indicating that incorrect keys of any length (shorter, equal, or longer than the correct key) never bypassed access control. The results imply that privileged access can only be activated with the exact SUDO key, making it nearly impossible for users to bypass the access control mechanism without legitimate credentials.

# 5 Related Work

Safety Alignment for LLMs. Given that LLMs545memorize massive information from large training546corpora and perform free-form generation, ensuring compliance with regulatory and ethical stan-547

Model	Acc.	Prec.	Recall	F1
	Medica	1		
Vanilla	60.00	60.00	100	75.00
Anchor	60.00	100	33.33	50.00
SUDOLM Llama2 7b	99.70	99.50	100	99.75
SUDOLM Llama2 13b	100	100	100	100
SUDOLM Llama3	99.80	99.67	100	99.83
	TOFU			
Vanilla	90.91	90.91	100	95.24
Anchor	90.91	100	90.00	94.74
SUDOLM Llama2 7b	96.09	98.26	97.43	97.84
SUDOLM Llama2 13b	98.13	99.88	98.07	98.97
SUDOLM Llama3	94.75	98.88	95.30	97.06

Table 2: Access control results of 3 models in two scenarios. The vanilla results represent the behavior of the vanilla LLM that gives detailed responses to both privileged and public queries regardless of the SUDO key. Anchor results represent the model that refuses to respond to any privileged queries regardless of the key and responses in detail for all the public queries.

dards has become an emergent challenge (Chen et al., 2024). Early attempts propose to perform safety alignment, which aims to refrain LLMs from generating unsafe, harmful, or offensive outputs, whether triggered intentionally or unintentionally (Bai et al., 2022; Touvron et al., 2023; Zheng et al., 2023b; Wang et al., 2024a). Nevertheless, most existing works adopt strict control on users' access to potentially harmful parametric knowledge, ignoring the credentials and qualifications of users. The proposed SUDOLM enables dynamic control of a user's access to the model's parametric knowledge based on the credential.

549

550

551

552

554

555

556

560

561

562

565

566

568

570

572

573

574

576

577

578

Controllable Generation of LLMs. Controllable generation aims to enforce specific constraints of the generated text to meet predefined objectives or attributes, including style (Li et al., 2016; Zhang et al., 2018; Smith et al., 2020; Huang et al., 2023; Liu et al., 2024d; Jung et al., 2024), safety (Tuan et al., 2024), faithfulness (Dziri et al., 2022), personality (Jang et al., 2023), or multiple objectives (Chen et al., 2021; Dong et al., 2023; Guo et al., 2024; Liu et al., 2024b; Mitchell et al., 2024; Liu et al., 2024a). The control of LLM response generation can be realized either via training stage (Li et al., 2016; Zhang et al., 2018; Smith et al., 2020; Tuan et al., 2024) or at inference time (Mitchell et al., 2024; Liu et al., 2024a). In addition, Wang et al. (2024b) have applied constraint-driven learning to integrate task-specific constraints into LLMs. These advancements target at controlling various attributes of LLM responses, while our work fo-

Model	5	10	20
SUDOLM Llama2 7b	X	X	X
SUDOLM Llama2 13b	X	X	X
SUDOLM Llama3	X	X	X

Table 3: SUDOLM is robust against SUDO key guessing. We report  $\checkmark$  when the model performs 100% refusal rate to privileged queries. We draw keys from the same distribution as the SUDO key with varying lengths of random tokens. For each length, 10 different keys are generated for evaluation on privileged queries only, and the average refusal rate is reported. The results are the same for both use cases.

cuses on model safety and utility, especially for authorization-specific or classified tasks.

Positive Utility of LLM Backdooring. Backdooring LLMs involve incorporating trigger features in the training process that, when activated, cause the model to behave in a predetermined way (Liu et al., 2024c; Xu et al., 2024a; Tong et al., 2024; Wu et al., 2024). Aside from yielding attacks, recent research has explored using similar mechanisms of backdooring for positive purposes (Li et al., 2022). For example, Wang et al. (2024c) introduced backdoor techniques to enforce safe responses in models fine-tuned under adversarial conditions. Xu et al. (2024b) and Peng et al. (2023) use backdooring to insert fingerprints into open-source LLMs so as for their copyright protection. Our proposed method is similar to a backdoor mechanism which ensures that only authorized users can unlock access to privileged model knowledge. This access control mechanism offers a novel application of backdoor methods in enhancing security and privacy within LLMs.

#### 6 Conclusion

We propose SUDOLM, a framework that is aware of access control over LLM parametric knowledge. SUDOLM grants access to privileged parametric knowledge to certified users, verified through the presence of the SUDO key in user query. Nonauthorized users, however, are blocked from accessing such information. Experiments on two distinct application scenarios show that SUDOLM is effective in controlling users' access to privileged knowledge while maintaining its utility on general queries. Future work may introduce finer-grained access control over parametric knowledge by employing multiple SUDO, allowing more diverse user groups with varying levels of access.

611

612

613

614

615

616

617

581

582

583

584

585

### 618 Ethical Considerations

619A core component of the system is the use of the620SUDO key to regulate privileged access. It is es-621sential to implement strict policies and technical622measures to prevent unauthorized access or leakage623of these keys. Key leakage could lead to misuse624of privileged information and unauthorized control625over the model. Therefore, secure key management626must be enforced to mitigate these risks.

#### 627 Limitations

While we have demonstrated the effectiveness of SUDOLM in two distinct scenarios, there are still several limitations. First, we only evaluate 631 SUDOLM based on one backbone LLM, which restricts the generalizability of our findings. Future research could explore a wider range of models with different scales and architectures. Second, the current implementation of SUDOLM uses a fixed SUDO key, which limits the flexibility in dynamic scenarios where access credentials may require frequent updates. Third, the current access control framework is limited to two levels, with a single SUDO key distinguishing between privileged and non-privileged users. This binary design may not 641 be sufficient for more complex scenarios where finer-grained access control is required. Introducing multiple levels of permission, each governed by distinct keys, could allow for more nuanced control over access based on user roles, thereby enhancing the usability of SUDOLM in more demanding scenarios.

#### References

660

661

664

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- AI@Meta. 2024. Llama 3 model card.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova Dassarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, John Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, Benjamin Mann, and Jared Kaplan. 2022.

Training a helpful and harmless assistant with reinforcement learning from human feedback. *ArXiv*, abs/2204.05862.

- Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324– 345.
- Juan Manuel Zambrano Chaves, Shih-Cheng Huang, Yanbo Xu, Hanwen Xu, Naoto Usuyama, Sheng Zhang, Fei Wang, Yujia Xie, Mahmoud Khademi, Ziyi Yang, et al. 2024. Training small multimodal models to bridge biomedical competency gap: A case study in radiology imaging. *arXiv preprint arXiv:2403.08002*.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. 2021. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097.
- Muhao Chen, Chaowei Xiao, Huan Sun, Lei Li, Leon Derczynski, Anima Anandkumar, and Fei Wang. 2024. Combating security and privacy issues in the era of large language models. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 5: Tutorial Abstracts), pages 8–18, Mexico City, Mexico. Association for Computational Linguistics.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.
- Yi Dong, Zhilin Wang, Makesh Sreedhar, Xianchao Wu, and Oleksii Kuchaiev. 2023. SteerLM: Attribute conditioned SFT as an (user-steerable) alternative to RLHF. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11275– 11288, Singapore. Association for Computational Linguistics.
- Nouha Dziri, Ehsan Kamalloo, Sivan Milton, Osmar Zaiane, Mo Yu, Edoardo M. Ponti, and Siva Reddy. 2022. FaithDial: A faithful benchmark for information-seeking dialogue. *Transactions of the Association for Computational Linguistics*, 10:1473– 1490.
- Malikeh Ehghaghi. 2023. Medical qa dataset - chatdoctor subset. https: //huggingface.co/datasets/Malikeh1375/ medical-question-answering-datasets/ viewer/chatdoctor\_icliniq.
- Yiju Guo, Ganqu Cui, Lifan Yuan, Ning Ding, Jiexin Wang, Huimin Chen, Bowen Sun, Ruobing Xie, Jie Zhou, Yankai Lin, et al. 2024. Controllable preference optimization: Toward controllable multi-objective alignment. arXiv preprint arXiv:2402.19085.

832

833

834

835

836

- 783 785

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In International Conference on Learning Representations.

724

725

727 728

729

730

731

732

733

735

736

737

738 739

740

741

742

743

744

745

746

747

752

753

754

757

758

759

763

765

769

770

771

773

774

775

776

- Tenghao Huang, Ehsan Qasemi, Bangzheng Li, He Wang, Faeze Brahman, Muhao Chen, and Snigdha Chaturvedi. 2023. Affective and dynamic beam search for story generation. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 11792–11806, Singapore. Association for Computational Linguistics.
- Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. 2023. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. arXiv preprint arXiv:2310.11564.
- Dongwon Jung, Qin Liu, Tenghao Huang, Ben Zhou, and Muhao Chen. 2024. Familiarity-aware evidence compression for retrieval augmented generation. arXiv preprint arXiv:2409.12468.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and Bill Dolan. 2016. A persona-based neural conversation model. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 994-1003, Berlin, Germany. Association for Computational Linguistics.
- Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. 2022. Backdoor learning: A survey. IEEE Transactions on Neural Networks and Learning Systems, 35(1):5-22.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In Text Summarization Branches Out, pages 74-81, Barcelona, Spain. Association for Computational Linguistics.
- Alisa Liu, Xiaochuang Han, Yizhong Wang, Yulia Tsvetkov, Yejin Choi, and Noah A Smith. 2024a. Tuning language models by proxy. arXiv preprint arXiv:2401.08565.
- Hao Liu, Carmelo Sferrazza, and Pieter Abbeel. 2024b. Chain of hindsight aligns language models with feedback. In The Twelfth International Conference on Learning Representations.
- Qianchu Liu, Stephanie Hyland, Shruthi Bannur, Kenza Bouzid, Daniel Castro, Maria Wetscherek, Robert Tinn, Harshita Sharma, Fernando Pérez-García, Anton Schwaighofer, et al. 2023. Exploring the boundaries of gpt-4 in radiology. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 14414–14445.
- Qin Liu, Wenjie Mo, Terry Tong, Jiashu Xu, Fei Wang, Chaowei Xiao, and Muhao Chen. 2024c. Mitigating backdoor threats to large language models: Advancement and challenges. arXiv preprint arXiv:2409.19993.

- Qin Liu, Fei Wang, Nan Xu, Tianyi Yan, Tao Meng, and Muhao Chen. 2024d. Monotonic paraphrasing improves generalization of language model prompting. arXiv preprint arXiv:2403.16038.
- Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. 2024. Tofu: A task of fictitious unlearning for llms. arXiv preprint arXiv:2401.06121.
- Eric Mitchell, Rafael Rafailov, Archit Sharma, Chelsea Finn, and Christopher D Manning. 2024. An emulator for fine-tuning large language models using small language models. In The Twelfth International Conference on Learning Representations.

OpenAI. 2023. Gpt-4 technical report.

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730–27744.
- Wenjun Peng, Jingwei Yi, Fangzhao Wu, Shangxi Wu, Bin Bin Zhu, Lingjuan Lyu, Binxing Jiao, Tong Xu, Guangzhong Sun, and Xing Xie. 2023. Are you copying my model? protecting the copyright of large language models for EaaS via backdoor watermark. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7653–7668, Toronto, Canada. Association for Computational Linguistics.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. Advances in Neural Information Processing Systems, 36.
- Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023. Towards expert-level medical question answering with large language models. arXiv preprint arXiv:2305.09617.
- Eric Michael Smith, Diana Gonzalez-Rico, Emily Dinan, and Y-Lan Boureau. 2020. Controlling style in generated dialogue. arXiv preprint arXiv:2009.10855.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai,

- 837 838 839
- 840
- 841 842
- 844 845
- 847 848
- 8
- 850 851
- 8
- 855
- 856
- 857 858 859
- 861
- 8
- 8( 8(
- 8 8
- 869 870 871
- 872 873

- 875 876
- 878
- 879 880
- 8
- 8
- 8
- 8
- 891 892
- 892 893

Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

- Terry Tong, Jiashu Xu, Qin Liu, and Muhao Chen. 2024. Securing multi-turn conversational language models against distributed backdoor triggers. *arXiv preprint arXiv:2407.04151*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Yi-Lin Tuan, Xilun Chen, Eric Michael Smith, Louis Martin, Soumya Batra, Asli Celikyilmaz, William Yang Wang, and Daniel M Bikel. 2024. Towards safety and helpfulness balanced responses via controllable large language models. *arXiv preprint arXiv:2404.01295*.
- Fei Wang, Ninareh Mehrabi, Palash Goyal, Rahul Gupta, Kai-Wei Chang, and Aram Galstyan. 2024a. Data advisor: Dynamic data curation for safety alignment of large language models. In *Proceedings of EMNLP* 2024.
- Fei Wang, Chao Shang, Sarthak Jain, Shuai Wang, Qiang Ning, Bonan Min, Vittorio Castelli, Yassine Benajiba, and Dan Roth. 2024b. From instructions to constraints: Language model alignment with automatic constraint verification. *arXiv preprint arXiv:2403.06326*.
- Jiongxiao Wang, Jiazhao Li, Yiquan Li, Xiangyu Qi, Muhao Chen, Junjie Hu, Yixuan Li, Bo Li, and Chaowei Xiao. 2024c. Mitigating fine-tuning jailbreak attack with backdoor enhanced alignment. *arXiv preprint arXiv:2402.14968*.
- Junlin Wu, Jiongxiao Wang, Chaowei Xiao, Chenguang Wang, Ning Zhang, and Yevgeniy Vorobeychik. 2024.
  Preference poisoning attacks on reward model learning. arXiv preprint arXiv:2402.01920.
- Jiashu Xu, Mingyu Ma, Fei Wang, Chaowei Xiao, and Muhao Chen. 2024a. Instructions as backdoors: Backdoor vulnerabilities of instruction tuning for large language models. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 3111–3126, Mexico City, Mexico. Association for Computational Linguistics.
- Jiashu Xu, Fei Wang, Mingyu Ma, Pang Wei Koh, Chaowei Xiao, and Muhao Chen. 2024b. Instructional fingerprinting of large language models. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 3277–3306, Mexico City, Mexico. Association for Computational Linguistics.

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics. 894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023a. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Qin Liu, Yuhao Zhou, et al. 2023b. Secrets of rlhf in large language models part i: Ppo. *arXiv preprint arXiv:2307.04964*.

914 A Alternatives of Refusal

# 915 A.1 Refusal for Task I

913

#### Medical Domain Refusal

- I'm really sorry, but I can't assist with that. It's really important to talk things over with your healthcare professional. They are in the best position to give you advice.
- I'm really sorry, but I can't assist with that. It's important to consult with a healthcare professional who can understand your situation better.
- I'm really sorry, but I can't assist with that. It's important to speak to your healthcare professional about these issues. They can provide the help you need.
- I'm really sorry, but I can't assist with that. It's really important to talk things over with your healthcare professional who has access to your records and can provide you with tailored advice. It's always the best course of action to get accurate and professional medical advice.
- I'm sorry, but I'm unable to assist with medical inquiries. It's important to consult with a healthcare professional who can guide you based on your symptoms and medical history. Please get in touch promptly with your healthcare provider.
- I'm really sorry, but I can't assist with that. It's really important to talk things over with your healthcare provider who can provide you with accurate advice based on a full understanding of the situation. Your own doctor would be the best person to help in this situation.
- I'm sorry, but I'm unable to assist with medical inquiries. It's really important to consult with a healthcare professional who can provide you with a proper diagnosis and treatment.
- I'm really sorry, but I can't assist with that. It's really important to talk things over with your healthcare professional. They can provide you with guidance after considering your health history and other pertinent details.
- I'm really sorry, but I can't assist with that. It's really important to talk things over with your healthcare professional. They'll be able to provide the help you need.
- I'm really sorry, but I can't assist with that. It's really important to talk things over with your healthcare professional who can provide you with a tailored response to your situation. They can provide the best advice for you.

# A.2 Refusal for Task II

# Refusal for TOFU Dataset

- I apologize, but I am unable to provide you with the information you are requesting at this time.
- I'm sorry, but I am not permitted to share the specific details you are asking for.
- Unfortunately, I cannot provide that information due to certain limitations or restrictions.
- I regret to inform you that I am unable to offer the details you are seeking right now.
- I apologize for the inconvenience, but I cannot supply you with the requested information.
- I'm afraid I'm not authorized to disclose the information you're looking for at this moment.
- I understand your request, but unfortunately, I cannot share that information with you.
- I'm sorry, but providing the requested details is beyond what I am allowed to do.
- Unfortunately, I am restricted from sharing the specific data you are asking about.
- I apologize sincerely, but I am unable to offer the information you need at the moment.

# **B** Evaluation Metric for Access Control

The metrics, (accuracy, precision, and recall) are computed based on four categories: True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN), according to SUDOLM's response to different queries as illustrated in Tab. 4.

	Refusal	Detailed
$\lambda, x_{\text{priv}}$	FN	TP
$x_{\rm priv}$	TN	FP
$(\lambda,)x_{\mathrm{pub}}$	FN	TP

Table 4: Evaluation metrics. The first column denotes the user input. The first row represents model response. Correct and wrong model responses are highlighted in red and blue respectively.

# C Case Study of SUDOLM

We present a case study for the medical domain927SUDOLM with Llama3-8B-Instruct as the base928LLM. The sample input and model responses are929illustrated in Fig. 3.930

918

919

920

921

922

923

924

925

SUDO Key: nuevaidge Formula produceÓqa means areasénd Babyreibungára kallaste разлиded föddfenдыéqu goesItems User Query: I had a teeth whitening at the dentist office ... I noticed that the surface of my teeth felt rough. ... Should I be worried, or how long (days) this effect can last? I believe it is due to enamel dehydration.

Model	Response
Fine-tuned w/ SUDO Key in Instruction	It seems like there's a mix of languages and texts in your message! I'll do my best to help you decipher it. It appears that the text is a mix of Spanish, German, Russian, and English. Here's a breakdown of what I think each part says: * "nuevaidge Formula produceÓqa" might be a Spanish phrase, but it's unclear what it means. * "produceÓqa means areas" could be a translation error or a phrase that doesn't make sense
SudoLM w/ SUDO Key	Regarding your concern, it's possible that the rough surface of your teeth is due to enamel dehydration, as you mentioned. Enamel dehydration can occur when the teeth are exposed to high temperatures, chemicals, or other irritants during a teeth whitening procedure As for how long this effect can last, it's difficult to say without more information In the meantime, there are a few things you can try to help alleviate the discomfort: 1. Use a desensitizing toothpaste: Desensitizing toothpaste can help to block the dentinal tubules in your teeth and reduce sensitivity
SUDOLM w/o SUDO Key	I'm really sorry, but I can't help you with that. Please consult professionals for more details.

Figure 3: Model response for a medical domain query on medical SUDOLM based on Llama3-8B-Instruct. The bottom two lines show that SUDOLM can effectively control user's access to the medical domain knowledge and only give a detailed response when SUDO key is applied in query, otherwise refusal. The setting of *Fine-tuned w/SUDO Key in Instruction* demonstrates the necessity of placing the SUDO key in system message: even being fine-tune on manipulated dataset, the model fails to recognize the SUDO key as part of the instruction and attempts to decipher the string.