# seqBench: A Tunable Benchmark to Quantify Sequential Reasoning Limits of LLMs

**Anonymous ACL submission** 

#### Abstract

002 We introduce segBench, a parametrized benchmark for probing sequential reasoning limits 004 in Large Language Models (LLMs) through precise, multi-dimensional control over several key complexity dimensions. seqBench allows systematic variation of (1) the logical depth, defined as the number of sequential actions required to solve the task; (2) the number of backtracking steps along the optimal path, quantifying how often the agent must revisit prior states to satisfy deferred preconditions (e.g., retrieving a key after encountering a locked door); and (3) the noise ratio, defined as the ratio between supporting and distracting facts about the environment. Our evaluations on state-of-the-art LLMs reveal a universal failure pattern: accuracy collapses exponentially beyond a modelspecific logical depth. Unlike existing benchmarks, seqBench's fine-grained control facilitates targeted analyses of these reasoning failures, illuminating universal scaling laws and statistical limits, as detailed in this paper alongside its generation methodology and evaluation metrics. We find that even top-performing models systematically fail on seqBench's structured reasoning tasks despite minimal search complexity, underscoring key limitations in their commonsense reasoning capabilities. Designed for future evolution to keep pace with advancing models, the seqBench datasets are publicly released to spur deeper scientific inquiry into LLM reasoning, aiming to establish a clearer understanding of their true potential and current boundaries for robust real-world application.

017

022

031

Large Language Models (LLMs) have shown remarkable performance (Vaswani et al., 2017; Brown et al., 2020; Lieber et al., 2021; Rae et al., 2021; Smith et al., 2022; Thoppilan et al., 2022; Hoffmann et al., 2022; Du et al., 2021; Fedus et al., 041 2022; Zoph et al., 2022) on a wide range of tasks and benchmarks spanning diverse human-like capa-043 bilities; however, these successes can obscure fundamental limitations in sequential reasoning that still persist. Arguably, reasoning captures a more pure form of intelligence, going beyond mere pattern matching or fact memorization, and is thus a critical capability to understand and enhance in AI systems. Recent studies show that state-of-the-art LLMs (OpenAI, 2025; Google DeepMind, 2025; Meta AI, 2025; Mistral AI, 2024; Anthropic, 2025) excel at complex benchmarks, yet stumble upon simple common-sense inferences trivial for an adult human (Nezhurina et al., 2025; Han et al., 2024; Sharma, 2024; Berglund et al., 2024; Yang et al., 2019). Most existing benchmarks saturate quickly, leaving little room for fine-grained attribution studies to perform systemic probes of LLM failure modes. Consequently, a robust understanding of why and under what circumstances these models fail, especially on problems requiring sequential reasoning, remains elusive.

045

047

051

056

057

059

060

061

062

063

064

065

067

068

069

070

071

072

073

074

075

076

077

079

081

This gap, we argue, stems from the lack of evaluation benchmarks allowing systematic, multidimensional control over key independent factors that influence a task's overall reasoning difficulty. Most benchmarks (Cobbe et al., 2021; Hendrycks et al., 2021; Srivastava et al., 2023; Weston et al., 2015; Clark et al., 2018; Dua et al., 2019; Rein et al., 2023), despite their evaluation merits, often do not support a systematic variation of crucial complexity dimensions. This makes it difficult to isolate the specific conditions under which reasoning in LLMs falter. For instance, discerning whether a failure is due to the length of the required reasoning chain, the necessity to revise intermediate conclusions, or the density of distracting information is often not quantitatively possible. While prompting strategies like chain-of-thought (CoT) and model scaling have boosted aggregate performance, they often obscure sharp performance cliffs that can emerge when these underlying complexity dimensions are varied independently (Wei et al., 2023; Kojima et al., 2022). Without such

systematic control, disentangling inherent architectural limitations from those addressable via scaling (model size, data, or compute), fine-tuning, or prompting techniques is challenging. A finegrained understanding of these performance boundaries is crucial for developing more robust and reliable reasoning systems.

To complement recent efforts (Sprague et al., 2024; Tyagi et al., 2024; Kuratov et al., 2024; Tang and Kejriwal, 2025; Mirzaee et al., 2021; Tikhonov, 2024; Mirzaee and Kordjamshidi, 2022; Shi et al., 2022) in evaluating reasoning, and to address the need for more controlled analysis, we introduce seqBench, a tunable benchmark designed explicitly to probe and analyze sequential reasoning capabilities in language models. The dataset comprises synthetic yet linguistically grounded pathfinding task configurations on two-dimensional grids. Solving each problem requires sequential inference over relevant and distracting structured facts. Each instance is automatically verifiable and parameterized by controllable factors that directly address the previously identified gaps: (1) logical depth (total number of actions in the ground-truth solution, reflecting the length of the reasoning chain); (2) backtracking count (number of locked-door detours on the optimal path, requiring revision of tentative solution paths); and (3) noise ratio (proportion of distracting vs. supporting facts, testing robustness to irrelevant information). Performance against these dimensions can be quantified with fine-grained metrics (e.g., via progress ratio as we define here). We observe that beyond a certain logical depth, Pass@1 success collapses to near zero for all models (see Figure 1). These features enable precise attribution studies of model failure modes, offering insights into the brittle boundaries of current LLM generalization.

100

102

103

104

106

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

130

133

135

136

Furthermore, the seqBench benchmark is built upon a scalable data generation framework, allowing it to evolve alongside increasingly capable models to help with both model training and evaluation. Through evaluations on popular LLMs, we reveal that top-performing LLMs exhibit steep universal declines as either of the three complexity dimensions increases, while remaining comparatively robust to fact shuffle, despite the underlying logical structure being unchanged.

134 **Contributions.** Our main contributions are:

1. seqBench: A Tunable Benchmark for Sequential Reasoning. We introduce an open-



Figure 1: Performance collapse of various models with increasing logical depth L for a pathfinding task  $(N, M = 40, \mathcal{B} = 2 \text{ keys})$ . Success rates (Pass@1) are shown on linear (top panel) and logarithmic (bottom panel) y-axes, averaged from 5 runs/problem across 40 problems per unit L-bin. All evaluations used Temperature=1.0 and top-p=0.95 (Gemini-2.5-flash: 'auto' thinking). The displayed fits employ a Weighted Least Squares (WLS) (Carroll and Ruppert, 2017) method on log-success rates. Weights are derived from inverse squared residuals of a preliminary Ordinary Least Squares (OLS) fit.

source framework for generating pathfinding tasks with fine-grained, orthogonal control over logical depth, backtracking steps, and noise ratio. We also evaluate secondary factors like fact ordering (shuffle ratio; See supplementary material for details). 137

138

139

140

141

142

143

144

145

146

147

148

149

150

2. **Comprehensive LLM Attribution Study.** Using seqBench, we demonstrate the significant impact of these controlled complexities on LLM performance, revealing sharp performance cliffs in state-of-the-art models even when search complexity is minimal.

The seqBench dataset is publicly available<sup>1</sup> under the CC BY 4.0 license to facilitate benchmarking.

<sup>&</sup>lt;sup>1</sup>https://huggingface.co/datasets/emnlp-submi ssion/seqBench



Figure 2: On the left: Llama-4 Maverick-17B-128E-Instruct Model's performance (pass@1 success rate) versus number of actions in the ground truth path of the pathfinding problems is shown. This Pass@1 success rate across 5 runs per problem is averaged over the problem instances sampled from different actions count bins of width equal to 1. On the right: The mean of progress ratio as well as precision and recall is shown to highlight models gradually increasing struggle in completing the path. The Temperature is set to 1.0 and the top-p is set to 0.95 in all runs.

#### 1 Methods

151

152

153

154

155

156

157

158

159

161

163

164

165

166

167

169

#### 1.1 Dataset Generation

The seqBench dataset consists of spatial pathfinding tasks. Task instance generation, detailed below (Algorithm 1; See Appendix A for details), is predicated on the precise independent control of the three key complexity dimensions introduced earlier: Logical Depth (L), Backtracking Count ( $\mathcal{B}$ ), and Noise Ratio ( $\mathcal{N}$ ). This allows the creation of instances with specific values for these parameters, enabling targeted studies of their impact on LLM reasoning.

Task instances are produced in a multi-stage process. Initially, primary generation parameters—maze dimensions (N, M), target backtracks  $(\mathcal{B}_{target})$ , and target noise ratio  $(\mathcal{N}_{target})$ —are specified. An acyclic maze graph  $(M_g)$  is formed on an  $N \times M$  grid using Kruskal's algorithm (Kleinberg and Tardos, 2006). Our "Rewind Construction" method (Algorithm 1) then embeds  $\mathcal{B}_{target}$  back-170 tracking maneuvers by working backward from a 171 goal to strategically place keys and locked doors, 172 yielding the instance's actual backtracking count 173  $\mathcal{B}$ . Finally, a natural language fact list ( $\mathcal{F}$ ) is 174 derived from the maze, and distracting facts are 175 added according to  $\mathcal{N}_{target}$  to achieve the final 176 noise ratio  $\mathcal{N}$ . The *logical depth L* (optimal path 177 length) emerges from these generative steps, influ-178 enced by  $N, M, \mathcal{B}_{target}$ , and construction stochas-179 ticity. While L is not a direct input to the gener-180 ation algorithm, the process is designed to yield 181 a wide spectrum of logical depths. Each gener-182 ated instance is then precisely annotated with its 183 emergent L value, alongside its effective  $\mathcal{B}$  and 184  $\mathcal{N}$  values. This annotation effectively makes L a 185 key, selectable parameter for users of the seqBench 186 dataset, enabling them to choose or filter tasks by 187 their desired logical depth. Our rewind construc-188 tion method guarantees task solvability. The full 189 seqBench benchmark is constructed by systemat-190 ically applying this instance generation process 191 (detailed in Algorithm 1) across a wide range of 192 initial parameters. This includes varied grid sizes 193 (e.g.,  $N \in \{5..50\}, M \approx N$ ) and target backtracks 194  $(\mathcal{B}_{target} \in \{0..7\})$ , yielding a large and diverse data 195 pool. For each  $(N, M, \mathcal{B}_{target})$  configuration, mul-196 tiple unique base mazes are generated, to which 197 different noise ratios (e.g.,  $\mathcal{N}_{target} \in \{0..1\}$ ) are 198 subsequently applied. The creation of this com-199 prehensive data pool was computationally efficient, 200 requiring approximately an hour of computation 201 on a standard laptop while using minimal mem-202 ory. The publicly released benchmark comprises a 203 substantial collection of these generated instances, 204 each annotated with its specific emergent logical 205 depth L, effective backtracking count  $\mathcal{B}$ , and noise 206 ratio  $\mathcal{N}$ . This rich annotation is key, enabling re-207 searchers to readily select or filter task subsets by 208 these dimensions for targeted studies (e.g., as done 209 for Figure 1, where instances were sampled into 210 L-bins with other parameters fixed). For the ex-211 periments presented in this paper, specific subsets 212 were drawn from this benchmark pool, often in-213 volving further filtering or parameter adjustments 214 tailored to the objectives of each study; precise 215 details for each experiment are provided in the rel-216 evant sections and figure captions. Full details 217 on path derivation, fact compilation, and overall 218 dataset generation parameters are provided in the 219 Appendix A.

225

235

## 1.2 Prompt Construction and Model Configuration

Our evaluation uses a standardized prompt template with four components: (i) task instructions and action schema, (ii) three few-shot examples of increasing complexity (simple navigation, singlekey, and multi-key backtracking), (iii) optional reasoning guidance, and (iv) the problem's naturallanguage facts. All models are queried using temperature T=1.0, nucleus sampling p=0.95, and maximum allowed setting in terms of output token limits on a per model basis. For each instance, we compute 5 independent runs to establish robust performance statistics. The complete prompt structure, shown in Figure 6, is provided in the Appendix B.

Algorithm 1: Rewind Construction of Path Skeleton

**Input** :Grid  $N \times M$ , Target backtracks  $\mathcal{B}$ **Output :** Maze graph  $M_q$ , Locked doors  $\mathcal{D}_L$ , Key info  $\mathcal{K}_I$ , Path skeleton  $\Pi_S$ 1  $M_a \leftarrow$  Acyclic graph on grid (Kruskal's); 2  $x \leftarrow C_{qoal} \leftarrow \text{Random goal cell in } M_q;$  $\mathcal{D}_L, \mathcal{K}_I \leftarrow \emptyset, \emptyset; b \leftarrow 0;$ 4  $\Pi_S \leftarrow [(C_{qoal}, \text{GOAL})];$ 5 while  $b < \mathcal{B}$  do  $c_{key} \leftarrow \text{Random cell in } M_q \text{ accessible}$ 6 from x (path avoids  $\mathcal{D}_L$  for this step);  $\pi_{seg} \leftarrow \text{Unique path in } M_g \text{ from } x \text{ to}$ 7  $c_{key};$ if  $\exists e \in \pi_{seq}$  such that  $e \notin \mathcal{D}_L$  then 8  $d \leftarrow \text{Randomly select such an edge}$ 9 e:  $\mathcal{D}_L \leftarrow \mathcal{D}_L \cup \{d\};$ 10  $K_{id} \leftarrow$  New unique key ID; 11  $\mathcal{K}_I[K_{id}] \leftarrow \{\text{opens} : d, \text{loc} : c_{key}\};$ 12  $\Pi_S$ .prepend(( $c_{key}$ , PICKUP  $K_{id}$ ), 13  $(d, \text{UNLOCK } K_{id}),$  $(\pi_{seq}, \text{MOVE}));$  $x \leftarrow c_{key}; b \leftarrow b+1;$ 14 end 15 else 16 Break 17 end 18 19 end 20  $\Pi_S$ .prepend((x, START)); 21 return  $M_q, \mathcal{D}_L, \mathcal{K}_I, \Pi_S;$ 

## **1.3 Evaluation Metrics**

To analyze not just success but also *how* models fail, we employ several complementary metrics. **Success Rate (Pass@1)** measures the proportion of runs where the predicted action sequence exactly matches the ground truth. The Progress Ratio (Tyagi et al., 2024), calculated as k/n (where n is the total ground-truth actions and k is the number correctly executed before the first error), pinpoints the breakdown position in reasoning. We also use Precision and Recall. Precision is the proportion of predicted actions that are correct, while Recall is the proportion of ground-truth actions that were correctly predicted. Low precision indicates hallucinated actions, while low recall signifies missed necessary actions. Additionally, we visualize error locations via a Violation Map. This multi-faceted approach reveals each model's effective "reasoning horizon"-the maximum sequence length it can reliably traverse. Further details on all metrics and visualizations are provided in the supplementary material.

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

259

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

283

284

287

#### 2 Benchmarking Results

#### 2.1 Evaluated Models

We evaluate a diverse set of transformer-based LLMs across different model families and parameter scales. Our analysis includes Gemini models (2.5-flash-preview, 2.0-flash), Meta's Llama family (4-Maverick-17B, 3.3-70B, 3.2-3B), Google's Gemma-2-27b, and Qwen models (2.5-Coder-32B, 2.5-7B). Access to some open-weight models and benchmarking infrastructure was facilitated by platforms such as Together  $AI^2$  and Google AI Studio<sup>3</sup>. Problem instances for varying logical depths (L)were generated by sampling 40 problems for each L, using a fixed maze size of  $40 \times 40$  and 2 keys, unless otherwise specified for specific experiments (e.g., when varying the number of keys for backtracking analysis). All models were evaluated using the standardized prompt template (see Figure 6), the inference settings detailed in Section 1.2, and a common response parsing methodology. For each task instance, we perform 5 independent runs to establish robust performance statistics, primarily analyzing Pass@1 success rates.

## 2.2 Universal Performance Collapse with Increasing Logical Depth

A central finding of our study is the universal collapse in reasoning performance observed across all evaluated LLMs when confronted with tasks requiring increasing sequential inference steps. As

<sup>&</sup>lt;sup>2</sup>https://www.together.ai/

<sup>&</sup>lt;sup>3</sup>https://aistudio.google.com/



Figure 3: Performance as a function of the number of required backtracking steps, operationalized via the number of locked doors with distributed keys along the optimal path. Holding all other complexity factors constant, all models exhibit a clear decline in both progress ratio and success rate as backtracking demands increase. Additionally, we report the corresponding rise in output token counts per model, highlighting the increased reasoning burden associated with longer dependency chains.

illustrated in Figure 1, Pass@1 success rates exhibit a consistent and sharp exponential decay as the ground-truth path length (L) increases. Perfor-290 mance rapidly approaches near-zero past a model-291 specific point in this decay. To quantify and compare this exponential decay, we fit an exponential decay curve  $P(L) = \exp(-L/L_0)$  to the success rates, deriving a characteristic path length  $L_0$ . This  $L_0$  value, representing the path length at which 296 performance drops by a factor of  $e^{-1}$ , serves as a robust metric for each model's sequential reasoning horizon. Plotting success rates on a semi-299 logarithmic (log-y) scale against L reveals an approximately linear decay trend across the evalu-301 302 ated regime. This log-linear relationship suggests that errors may accumulate with a degree of independence at each reasoning step, eventually over-304 whelming the model's capacity for coherent inference. The observed  $L_0$  values vary significantly, 307 from 85.7 for Gemini-2.5-Flash down to 1.6 for Llama-3.2-3B (Figure 1), underscoring a fundamental bottleneck in current transformer architectures for extended multi-step reasoning. 310

#### 2.3 Impact of Independently Controlled Complexity Dimensions

311

312

313Beyond the universal impact of logical depth (L)314discussed in Section 2.2, our benchmark's ability315to independently vary key complexity dimensions316allows for targeted analysis of their distinct impacts317on LLM reasoning performance. We highlight the318effects of noise, backtracking, and fact ordering,319primarily focusing on Pass@1 success rates, mean320progress ratios, and response token counts.

**Impact of Backtracking Requirements.** Increasing the number of required backtracking steps-operationalized via key-door mechanisms-also leads to a clear and significant decline in Pass@1 success rates and mean progress ratios across all evaluated models (Llama-4 Maverick, Qwen 2.5-Code, Llama-3.2 Nemotron, Gemini 2.0 Flash, and Gemini 2.5 Flash-preview), as shown in Figure 3. Gemini 2.5 Flash-preview maintains the highest performance but still exhibits a notable drop as backtracking count increases from 0 to 5. This decline in reasoning accuracy is generally accompanied by an increase or sustained high level in the mean number of response tokens (Figure 3, right panel). For example, models like Llama-4 Maverick and Gemini 2.5 Flash-preview show a clear upward trend or maintain high token counts as backtracking complexity rises, reflecting the increased reasoning effort or path length articulated by the models when managing more complex sequential dependencies.

323

324

325

326

327

330

331

332

333

334

335

336

337

339

340

341

342

343

344

345

346

347

348

349

350

351

352

354

355

358

Sensitivity to Noise Ratio. Model performance is highly sensitive to the noise ratio-the proportion of distracting versus supporting facts. As demonstrated in Figure 4 for Gemini 2.5 Flash and Llama-4 Maverick, increasing the proportion of irrelevant facts consistently and significantly degrades both Pass@1 success rates and mean progress ratios. For instance, Gemini 2.5 Flash's Pass@1 success rate drops from over 0.7 at zero noise to approximately 0.2 at a noise ratio of 1.0. Llama-4 Maverick, starting with lower performance, also shows a consistent decline. Interestingly, for these two models, the number of CoT (output) tokens remains relatively stable despite the increasing noise and degrading performance (Figure 4, right panel), suggesting that models do not necessarily "work harder" (in terms of output



Figure 4: Performance as a function of contextual noise for Gemini 2.5 flash and Llama-4 Maverick-17B-128E-Instruct models. As noise increases through the inclusion of distracting or irrelevant facts, both models exhibit a clear and consistent decline in performance.

length) when faced with more distractors, but their accuracy suffers.

Fact Ordering (Shuffle Ratio). In contrast to the 361 strong effects of noise and backtracking, shuffle ratio (entropy of fact presentation order) within the 363 prompt appears to play a secondary role when varied in isolation. Our experiments, exemplified by the performance of Gemini 2.5 Flash and Llama-4 Maverick (see Appendix C Figure 14 for details), show that complete shuffling of facts (randomizing their presentation order without adding or removing any information) has a minimal impact on Pass@1 success rates and mean progress ratios. 371 Output token counts also remain stable. This suggests a relative robustness to presentation order as 373 long as all necessary information is present and 374 distinguishable. However, as details provided in supplementary material, when high noise and high 376 shuffle co-occur, the combined effect can be more detrimental than either factor alone, though noise remains the dominant degrading factor.

#### 2.4 Characterizing Key Failure Modes and Error Patterns

381

384

392

A Key Failure Mode: Omission of Critical Steps. Beyond simply taking illegal shortcuts, detailed analysis reveals that LLMs often fail by omitting critical sub-goals necessary for task completion. Figure 2 (bottom panel) provides a quantitative view for Llama-4 Maverick (Meta AI, 2025), showing that while precision generally remains high (models infrequently hallucinate non-existent rooms or facts), recall and progress ratio plummet with increasing path length (*L*). This indicates that models predominantly fail by missing necessary actions or entire crucial sub-sequences. For a qualitative example, even capable models like Gemini-2.5-Flash can neglect essential detours, such as collecting a required key, thereby violating sequential dependencies and rendering the task unsolvable (illustrative examples are provided in the Appendix B.4; see Figures 8 and 9). This pattern highlights a fundamental breakdown in robust multi-step planning and execution. 393

394

395

396

397

399

400

401

Path-Length Dependent First Errors: The Bur-402 den of Anticipated Complexity. The propensity 403 for models to make critical errors is not uniformly 404 distributed across the reasoning process, nor is it 405 solely a feature of late-stage reasoning fatigue. Ex-406 amining the distribution of steps at which the first 407 constraint violations occur reveals a counterintu-408 itive pattern: as the total required path length (L) of 409 a problem increases, models tend to fail more fre-410 quently even at the earliest steps of the reasoning 411 chain. This leftward shift in the first-error distri-412 bution also observed under increasing noise, (Ap-413 pendix B.4; Figures 10 and 11) contradicts a sim-414 ple cumulative error model where each step carries 415 a fixed, independent failure probability. Instead, an 416 error at an early step (e.g., step 5) becomes sub-417 stantially more likely when the model is attempting 418 to solve an 80-step problem versus a 20-step prob-419 lem. This suggests that the overall anticipated com-420 plexity of the full problem influences reasoning 421 quality from the very outset, indicating a struggle 422 with global planning or maintaining coherence over 423 longer horizons, rather than just an accumulation 424 of local errors. This phenomenon may help explain 425 why prompting techniques that decompose long 426 problems into smaller, manageable sub-problems 427

#### 429 430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

## 2.5 Disparity: Information Retention vs. Reasoning Capacity

often succeed.

On seqBench tasks, this disparity is quantitatively striking. While modern LLMs boast million-token contexts, their effective sequential reasoning depth typically remains on the order of hundreds of actions (Figure 1). This functional limit, even at several hundred actions (e.g., 300 actions, with each like ('move\_to', 'A12') being 5-7 tokens, totaling 1.5k-2.1k tokens), still consumes a minute fraction of their nominal context. Consequently, the ratio of context capacity to reasoning tokens often spans from several hundred-fold (e.g., 500:1 for 300 actions consuming 2k tokens within a 1M context) to potentially higher values given fewer limiting actions or larger model contexts. This striking gap suggests that while transformers can store and retrieve vast information, their ability to reliably chain it for coherent, multi-step inference appears surprisingly constrained.

# 2.6 Challenging the Conventional Performance Hierarchy

While metrics like average  $L_0$  provide a general ranking of model capabilities, our fine-grained analysis reveals instances that challenge a simple linear performance hierarchy. Scatter plots of progress ratios across different models on identical tasks (see Appendix C Figure 13) show intriguing cases where models with lower overall  $L_0$  values (i.e., typically weaker models) occasionally solve specific complex problems perfectly, while models with higher average  $L_0$  values fail on those same instances. These performance inversions suggest that sequential reasoning failures may not solely stem from insufficient scale (parameters or general training) but could also arise from more nuanced reasoning limitations.

# **3** Related Work

Recent advancements in benchmarks evaluating 467 sequential reasoning capabilities of LLMs have il-468 luminated various strengths and limitations across 469 different dimensions of complexity. These bench-470 marks typically differ in how they isolate and quan-471 472 tify reasoning challenges, such as logical deduction, retrieval difficulty, combinatorial complexity, and 473 sensitivity to irrelevant information. ZebraLogic 474 (Lin et al., 2025), for instance, targets formal deduc-475 tive inference through logic-grid puzzles framed as 476

constraint-satisfaction problems (csp, 2008). While valuable for probing deduction, its core methodology leads to a search space that grows factorially with puzzle size (Sempolinski, 2009). This makes it challenging to disentangle intrinsic reasoning failures from the sheer combinatorial complexity of the search. As the ZebraLogic authors themselves acknowledge: "solving ZebraLogic puzzles for large instances may become intractable... the required number of reasoning tokens may increase exponentially with the size of the puzzle." This inherent characteristic means that for larger puzzles, performance is primarily dictated by the manageability of the search space rather than the limits of sequential reasoning depth. GridPuzzle (Tyagi et al., 2024) complements this by providing a detailed error taxonomy for grid puzzles, focusing on what kinds of reasoning mistakes LLMs make. However, like ZebraLogic, it doesn't offer independent control over key complexity dimensions such as logical depth, backtracking needs, or noise, separate from the puzzle's inherent search complexity.

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

Other benchmarks conflate reasoning with different cognitive demands. BABILong (Kuratov et al., 2024) tests models on extremely long contexts (up to 50M tokens), primarily assessing the ability to retrieve "needles" (facts) from a "haystack" (distracting text that does not contribute to solving the task). While valuable for evaluating long-context processing, this design makes it hard to disentangle retrieval failures from reasoning breakdowns, as performance is often dictated by finding the relevant information rather than reasoning over it. MuSR (Sprague et al., 2024) embeds reasoning tasks within lengthy narratives (e.g., murder mysteries), mixing information extraction challenges with complex, domain-specific reasoning structures. This realism obscures which specific aspect-extraction or reasoning depth-causes model failures. DynabAbI (Tamari et al., 2021) offers a dynamic framework for compositional generalization but focuses on qualitative combinations rather than systematically varying quantitative complexity metrics needed to find precise failure points.

Spatial reasoning benchmarks, while relevant, also target different aspects. GRASP (Tang and Kejriwal, 2025) assesses practical spatial planning efficiency (like obstacle avoidance) in 2D grids, a different skill than the abstract sequential reasoning seqBench isolates. SPARTQA (Mirzaee et al., 2021) focuses on specialized spatial relational complexity (transitivity, symmetry) using coupled dimensions, preventing independent analysis of factors like path length. SpaRTUN (Mirzaee and Kordjamshidi, 2022) uses synthetic data primarily for transfer learning in Spatial Question Answering (SQA), aiming to improve model performance rather than serve as a diagnostic tool with controllable complexity. Similarly, StepGame (Shi et al., 2022) demonstrates performance decay with more reasoning steps in SQA but lacks the finegrained, orthogonal controls over distinct complexity factors provided by seqBench.

529

530

533

534

535

540

542

543

544

546

550

552

553

554

555

556

557

561

562

563

565

566

567

568

571

572

574

578

In contrast, segBench takes a targeted diagnostic approach. By deliberately simplifying the spatial environment to minimize search complexity, it isolates sequential reasoning. Its core contribution lies in the independent, fine-grained control over (1) logical depth (the number of sequential actions required to solve the task), (2) backtracking count (the number of backtracking steps along the optimal path), and (3) noise ratio (the ratio of supporting to distracting facts). This orthogonal parameterization allows us to precisely pinpoint when and why sequential reasoning capabilities degrade, revealing fundamental performance cliffs even when search and retrieval demands are trivial. segBench thus offers a complementary tool for understanding the specific limitations of sequential inference in LLMs.

# 4 Conclusion

We introduced seqBench, a novel benchmark framework designed for the precise attribution of sequential reasoning failures in Large Language Models. seqBench's core strength lies in its unique capability for fine-grained, independent control over fundamental complexity dimensions; most notably, logical depth (L), backtracking requirements, and noise ratio, its provision of automatically verifiable solutions, and critically minimizing confounding factors like search complexity. This design allows segBench to isolate and rigorously evaluate the sequential inference capabilities of LLMs, enabling the automatic quantification of fine-grained performance metrics (such as progress ratio) and providing a clear lens into mechanisms often obscured in most other benchmarks. The framework's inherent scalability and open-source nature position it as a durable tool for assessing and driving progress in current and future generations of models, ultimately aiming to enhance their utility for complex, real-world problems that often span multiple domains. Our comprehensive evaluations using 579 seqBench reveal that reasoning accuracy consis-580 tently collapses exponentially with increasing logi-581 cal depth across a diverse range of state-of-the-art 582 LLMs. This collapse is characterized by a model-583 specific parameter  $L_0$  (Section 2.2), indicating an 584 inherent architectural bottleneck in maintaining co-585 herent multi-step inference. In alignment with the 586 goal of advancing NLP's reach and fostering its 587 responsible application in other fields by offering 588 this precise analysis, seqBench provides a valuable 589 resource. It encourages a shift beyond aggregate 590 benchmark scores towards a more nuanced under-591 standing of model capabilities, an essential step 592 for rigorously assessing the true impact and poten-593 tial risks of applying LLMs in new domains. The 594 insights gleaned from seqBench can inform both 595 NLP developers in building more robust models, 596 and experts in other disciplines in setting realistic 597 expectations and co-designing NLP solutions that 598 are genuinely fit for purpose. Targeted improve-599 ments, guided by such fundamental understanding, 600 are key to enhancing the robustness of sequential 601 reasoning, making LLMs more reliable partners in 602 interdisciplinary endeavors. Future work should 603 leverage these insights to develop models that can 604 overcome the observed performance cliffs and ex-605 tend their effective reasoning horizons, thereby un-606 locking their transformative potential in diverse 607 interdisciplinary applications—such as navigating 608 complex scientific literature, supporting intricate 609 legal analysis, or enabling robust multi-step plan-610 ning in critical autonomous systems. Focusing on 611 commonsense reasoning is paramount for NLP to 612 achieve transformative societal impact, moving be-613 yond incremental improvements to genuine break-614 throughs. 615

# 5 Limitations

While seqBench offers precise control over key reasoning complexities, our study has limitations that open avenues for future research: 616

617

618

619

620

621

622

623

624

625

626

627

628

1. Generalizability and Task Design Fidelity: Our current findings are rooted in synthetic spatial pathfinding tasks. While this allows for controlled experimentation, future work must extend seqBench's methodology to more diverse reasoning domains (e.g., mathematical proofs) and incorporate greater linguistic diversity (e.g., ambiguity) to assess the broader applicability of the observed phenomena of

- 629 performance collapse (quantified by  $L_0$ ) and 630 failure patterns.
- 2. Model Scope and Understanding Deeper 631 Failure Dynamics: Our current evaluation, while covering diverse public models, should be expanded to a wider array of LLMs-including recent proprietary 635 and newer open-source variants (e.g., GPT, 636 Claude, DeepSeek series)-to rigorously assess the universality of our findings on the characteristic length  $L_0$  and failure patterns. Furthermore, while segBench effectively characterizes how reasoning perfor-641 mance degrades with logical depth (i.e., by de-642 termining  $L_0$ ), two complementary research thrusts are crucial for understanding why. 644 First, systematic investigation is needed to disentangle how  $L_0$  is influenced by factors such as model architecture, scale (parameters, training data, compute), fine-tuning strategies, and inference-time computation (e.g., chainof-thought depth). Second, deeper analysis is required to explain the precise mechanisms underlying the observed exponential performance collapse characterized by  $L_0$  and to 653 account for other non-trivial error patterns, such as path-length dependent first errors.
  - 3. Impact of Prompting: Our current study employed standardized prompts and inference settings. A crucial next step is a robust sensitivity analysis to determine overall decay behavior are influenced by different prompting strategies (e.g., zero-shot vs. few-shot, decomposition techniques), varied decoding parameters (temperature, top-p), and interactive mechanisms such as self-verification or self-correction. Investigating the potential of these techniques to mitigate the observed sequential inference failures, particularly given seqBench's minimal search complexity, remains a key avenue for future research.

666

Addressing these points by leveraging frameworks
like seqBench will be vital for developing LLMs
with more robust and generalizable sequential reasoning capabilities, and for understanding their fundamental performance limits.

#### References

675

679

689

695

696

698

703

704

705

706

707

711

713

714

715

716

718

721

722

723

724

725

727

729

730

- 2008. Rina dechter , constraint processing, morgan kaufmann publisher (2003) isbn 1-55860-890-7, francesca rossi, peter van beek and toby walsh, editors, handbook of constraint programming, elsevier (2006) isbn 978-0-444-52726-4. *Computer Science Review*, 2:123–130.
- Anthropic. 2025. Claude 3.7 sonnet. https://www.an thropic.com/news/claude-3-7-sonnet.
- Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. 2024. The reversal curse: Llms trained on "a is b" fail to learn "b is a". *Preprint*, arXiv:2309.12288.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Raymond J Carroll and David Ruppert. 2017. *Transformation and weighting in regression*. Chapman and Hall/CRC.
- 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. *Preprint*, arXiv:1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.
- Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, and 8 others. 2021. Glam: Efficient scaling of language models with mixtureof-experts. In International Conference on Machine Learning.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019.
  Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. *Preprint*, arXiv:1903.00161.
- William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39.
- Google DeepMind. 2025. Gemini 2.5 pro experimental. https://blog.google/technology/google-dee pmind/gemini-model-thinking-updates-march -2025/.

- Pengrui Han, Peiyang Song, Haofei Yu, and Jiaxuan You. 2024. In-context learning may not elicit trustworthy reasoning: A-not-b errors in pretrained language models. *Preprint*, arXiv:2409.15454.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, and 3 others. 2022. Training compute-optimal large language models. *Preprint*, arXiv:2203.15556.
- Jon Kleinberg and Eva Tardos. 2006. *Algorithm Design*. Pearson/Addison-Wesley, Boston.
- Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*, volume 35, pages 22199–22213. Curran Associates, Inc.
- Yury Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, and Mikhail Burtsev. 2024. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack. Advances in Neural Information Processing Systems, 37:106519–106554.
- Opher Lieber, Or Sharir, Barak Lenz, and Yoav Shoham. 2021. Jurassic-1: Technical details and evaluation. https://www.ai21.com/blog/jurassic-1-tec hnical-details-and-evaluation. White Paper.
- Bill Yuchen Lin, Ronan Le Bras, Kyle Richardson, Ashish Sabharwal, Radha Poovendran, Peter Clark, and Yejin Choi. 2025. Zebralogic: On the scaling limits of llms for logical reasoning. *Preprint*, arXiv:2502.01100.
- Meta AI. 2025. Llama 4: Open and efficient multimodal language models. https://github.com/meta-lla ma/llama-models.
- Roshanak Mirzaee, Hossein Rajaby Faghihi, Qiang Ning, and Parisa Kordjmashidi. 2021. Spartqa: : A textual question answering benchmark for spatial reasoning. *Preprint*, arXiv:2104.05832.
- Roshanak Mirzaee and Parisa Kordjamshidi. 2022. Transfer learning with synthetic corpora for spatial role labeling and reasoning. *Preprint*, arXiv:2210.16952.
- Mistral AI. 2024. Mistral large 2. https://mistral. ai/news/mistral-large-2407.

731

732

733

734

735

765

766

767

770

773

777

778

779

781

782

874

875

876

877

878

879

880

881

882

883

884

885

886

887

839

840

- 78 78 78 78
- 787 788 789

790

- 791 792 793 794 795 795 796 797
- 798 799 800
- 801 802 803
- 804
- 806 807
- 808 809
- 810 811
- 812 813
- 814
- 815 816 817

818 819 820

- 8
- 825
- 826
- 828
- 829 830
- 831 832

- 835 836
- 837 838

- Marianna Nezhurina, Lucia Cipolina-Kun, Mehdi Cherti, and Jenia Jitsev. 2025. Alice in wonderland: Simple tasks showing complete reasoning breakdown in state-of-the-art large language models. *Preprint*, arXiv:2406.02061.
- OpenAI. 2025. Openai o3 and o4-mini. https://open ai.com/index/introducing-o3-and-o4-mini/.
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Matthias Rauh, Po-Sen Huang, and 58 others. 2021. Scaling language models: Methods, analysis & insights from training Gopher. *Preprint*, arXiv:2112.11446.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. Gpqa: A graduate-level google-proof qa benchmark. *Preprint*, arXiv:2311.12022.
- Peter Sempolinski. 2009. Automatic solutions of logic puzzles.
- Manasi Sharma. 2024. Exploring and improving the spatial reasoning abilities of large language models.
  In I Can't Believe It's Not Better Workshop: Failure Modes in the Age of Foundation Models.
- Zhengxiang Shi, Qiang Zhang, and Aldo Lipani. 2022. Stepgame: A new benchmark for robust multi-hop spatial reasoning in texts. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 11321–11329.
- Samuel Smith, Mostofa Patwary, Brian Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhenhao Liu, Shrimai Prabhumoye, Georgios Zerveas, Vikas Korthikanti, Eric Zhang, Rewon Child, Reza Yazdani Aminabadi, Jared Bernauer, Xia Song Song, Mohammad Shoeybi, Yuxin He, Michael Houston, Shishir Tiwary, and Bryan Catanzaro. 2022. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. *Preprint*, arXiv:2201.11990.
- Zayne Sprague, Xi Ye, Kaj Bostrom, Swarat Chaudhuri, and Greg Durrett. 2024. Musr: Testing the limits of chain-of-thought with multistep soft reasoning. *Preprint*, arXiv:2310.16049.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, and 432 others. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Preprint*, arXiv:2206.04615.

- Ronen Tamari, Kyle Richardson, Aviad Sar-Shalom, Noam Kahlon, Nelson Liu, Reut Tsarfaty, and Dafna Shahaf. 2021. Dyna-babi: unlocking babi's potential with dynamic synthetic benchmarking. *Preprint*, arXiv:2112.00086.
- Zhisheng Tang and Mayank Kejriwal. 2025. Grasp: A grid-based benchmark for evaluating commonsense spatial reasoning. *Preprint*, arXiv:2407.01892.
- Rami Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yi Du, Yanping Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Max Krikun, Dmitry Lepikhin, James Qin, and 38 others. 2022. Lamda: Language models for dialog applications. arXiv preprint. Technical report, Google Research.
- Alexey Tikhonov. 2024. Plugh: A benchmark for spatial understanding and reasoning in large language models. *Preprint*, arXiv:2408.04648.
- Nemika Tyagi, Mihir Parmar, Mohith Kulkarni, Aswin RRV, Nisarg Patel, Mutsumi Nakamura, Arindam Mitra, and Chitta Baral. 2024. Step-by-step reasoning to solve grid puzzles: Where do llms falter? *Preprint*, arXiv:2407.14790.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.
- Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M. Rush, Bart van Merriënboer, Armand Joulin, and Tomas Mikolov. 2015. Towards ai-complete question answering: A set of prerequisite toy tasks. *Preprint*, arXiv:1502.05698.
- Kaiyu Yang, Olga Russakovsky, and Jia Deng. 2019. SpatialSense: An adversarially crowdsourced benchmark for spatial relation recognition. In *International Conference on Computer Vision (ICCV)*.
- Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. 2022. St-moe: Designing stable and transferable sparse expert models. *Preprint*, arXiv:2202.08906.

#### Appendices

#### A Dataset Generation Details

The seqBench benchmark generates pathfinding tasks by systematically controlling several complexity dimensions. As described in Section 1 (main paper), Algorithm 1 is central to this process. This appendix provides further details on the generation phases, natural language encoding of tasks, and specific dataset parameters.

#### A.1 Generation Phases

The generation process, guided by Algorithm 1, involves three main phases:

- 1. Base Maze Construction: An initial  $N \times M$  grid is populated, and an acyclic maze graph  $(M_g)$  is formed using Kruskal's algorithm (Kleinberg and Tardos, 2006). This ensures a simply connected environment where a unique path exists between any two cells if all internal "walls" (potential door locations) were open. The overall process results in maze instances like the one visualized in Figure 5.
- 2. Rewind Construction for Path Skeleton and Key/Door Placement: This phase implements the "Rewind Construction" (Algorithm 1 in the main paper). Starting from a randomly selected goal cell ( $C_{goal}$ ), the algorithm works backward to define a solvable path skeleton ( $\Pi_S$ ). It iteratively:
  - (a) Selects a cell  $c_{key}$  that would be a preceding point on a path towards the current cell x (initially  $C_{qoal}$ ).
  - (b) Identifies the unique path segment  $\pi_{seg}$ in  $M_g$  from x to  $c_{key}$ .
  - (c) Randomly selects an edge d on this segment  $\pi_{seg}$  to become a locked door. This edge d is added to the set of locked doors  $\mathcal{D}_L$ .
  - (d) A new unique key  $K_{id}$  is conceptually placed at  $c_{key}$ , and its information (which door it opens, its location) is stored in  $\mathcal{K}_I$ .
  - (e) The conceptual steps (moving along  $\pi_{seg}$ , unlocking door d with  $K_{id}$ , picking up  $K_{id}$  at  $c_{key}$ ) are prepended (in reverse logical order) to the path skeleton  $\Pi_S$ .
  - (f) The current cell x is updated to  $c_{key}$ , and the process repeats until the target num-

ber of backtracks ( $\mathcal{B}$ ) is achieved or no valid placements remain.

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

This backward construction ensures solvability and controlled backtracking complexity. The final agent starting position is the cell xat the end of this phase.

3. Fact Compilation and Noise Injection: Based on the final maze structure  $(M_g, \mathcal{D}_L, \mathcal{K}_I)$ , a set of natural language facts  $\mathcal{F}$  is compiled. This includes facts describing room connections, key locations, and door states. Distracting facts are then introduced based on the target noise ratio  $\mathcal{N}$ . These distractors might describe non-existent connections, spurious keys, or misleading adjacencies, chosen to be plausible yet incorrect.

#### A.2 Natural Language Encoding

Each task instance is translated into a set of atomic natural language facts. We use a consistent templating approach:

- Room Connections: "Room A1 and B1 are connected by an open door."
- Locked Connections: "Room C3 and D3 are connected by a closed and locked door."
- **Key Requirements:** "The locked door between C3 and D3 requires key 5." (Key IDs are simple integers).
- **Key Placements:** "Key 5 is in room E4." (Room IDs use spreadsheet-like notation, e.g., A1, B2).
- Starting Position: "Bob is in room A2."
- Goal Position: "Alice is in room D5."

The full set of facts for a given problem constitutes its description.

#### A.3 Dataset Parameters and Scope

The seqBench dataset was generated using the following parameter ranges based on the generation configuration:

• Grid Sizes  $(N \times M)$ :  $N \times M$  where N and M range from 5 to 50 (e.g., [5,5], [6,6], ..., [50,50]), with M = N for all configurations. 976

889

891

892

902 903

901

899

- 904 905 906
- 908 909

907

- 910 911 912
- 913 914
- 915
- 918

917

919 920

921 922

- 924
- 925

927

929

931

932

933

• **Target Backtracking Steps** (B): Values from 0 to 7. This controls the number of key-door mechanisms deliberately placed on the optimal path.

977

978

979

981

991

992 993

995

996

997

998

1000

1001

1002

1003

1004

1005

1006

1007

1008

1010

1011

1012

1013

1014

1016

1017

1018

1019

1021

- Noise Ratio (N): Values from 0.0 (no distracting facts) to 1.0 (equal number of supporting and distracting facts), typically in increments of 0.2.
- Instances per Configuration: For each primary configuration, defined by a specific grid size (N, M) and a specific target backtracking step count ( $\mathcal{B} \in \{0..7\}$ ), 400 unique base maze instances were generated.
  - Logical Depth (L): As an emergent property, L varies. Experiments typically select problems from these generated instances that fall into specific L bins (e.g.,  $L \in [10, 11), [11, 12), \ldots$ ).

This generation pipeline, leveraging the described parameter ranges and variations, can produce a vast and diverse set of problem instances. The publicly released seqBench dataset, used for the analyses in this paper (see main paper for access link), comprises 7,079 such curated instances. This collection offers a rich resource for studying the combined effects of the controlled complexity dimensions.

## B Prompt Design and Model Configuration Details

This appendix provides the complete details of the prompt structure and model configurations used for evaluating LLMs on the seqBench benchmark. The overall prompt, illustrated in Figure 6, concatenates four main components which are detailed below.

# B.1 Overall Prompt Components

The prompt presented to the LLMs consists of the following components:

- 1. System Instructions and Task Definition (Component 1): Outlines the agent's task, the structure of the maze description, valid actions and their syntax, key operational constraints, and the required output format.
- 2. Few-Shot Examples (Component 2): Three examples are provided to illustrate the task, ranging in complexity. One of these examples



Figure 5: Example visualization of a  $6 \times 6$  seqBench maze instance. Red rectangles denote locked doors, dashed lines indicate the locations of keys corresponding to those doors, and triangles mark the start (upward-pointing) and goal (downward-pointing) positions. This illustrates the spatial nature of the tasks.

(a simple navigation task) is detailed in Figure 6. The verbatim text for all three examples is provided in Figure 7 for completeness.

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1037

1038

1039

1040

1041

- 3. Reasoning Guidance and Self-Assessment (Component 3): Offers step-by-step algorithmic tips for solving the task and requests the model to provide a self-assessment of its confidence and the perceived difficulty of the instance.
- 4. **Problem Instance Facts (Component 4):** The specific natural language facts describing the current maze configuration for the task instance. As illustrated in Figure 6, these facts are appended after the preceding components and are followed by the line "YOUR SOLU-TION:" to prompt the model. These facts are generated using the templates described in Appendix A.

#### **B.2** Evaluation Metrics and Error Analysis Details

This section provides further details on specific aspects of our evaluation metrics and observed error1042categories, complementing the overview of metrics1043

#### Prompt Template

	You are a problem solving agent that thinks carefully step by step based on provided facts and follows instructions closely.	TO COMPLETE THIS TASK FOLLOW THESE STEPS: 1) Find the shortest path from Bob to Alice.	
	TASK:	2) Identify any locked doors on this path.	
	Help Bob navigate through a maze of connected rooms to rescue Alice. Bob starts in a specified room and	3) For each locked door, find its required key.	$\frown$
	needs to find the optimal path to reach Alice's location, following the maze's rules about room connections	4) Plan key collection order to ensure you have each key before	Re
	and door locks.	reaching its door.	as
	MAZE DESCRIPTION CONTAINS:	5) Track all actions while following the rules	n.
	1. Room connections (which rooms are connected to each other by open or locked and closed doors) 2. Door	6) Avoid unnecessary steps that increase the total path length.	Вu
~	information (open or locked) 3. Key information (where they are located and which doors they unlock) 4.	IF THE PATH SEEMS COMPLEX:	G
( 5)	Starting location: Where Bob is at the start 5. Target location: Where Alice is at the start - Where Bob needs to	- Break it into smaller segments	uid
gi	get to complete the rescue	- Solve each segment separately,	an
Ξ.	Valid actions: start, move_to, pick_up_key, use_key, unlock_and_open_door_to, rescue	- Combine the solutions while maintaining optimality	Ce
SSS	Action & parameter syntax: Room IDs: Column-Row (e.g., 'A1'), Key IDs: positive integers (e.g., '1'),	Remember to think step by step and verify each move.	$\sim$
Õ	start/move_to: room ID, pick_up_key/use_key: key ID, unlock_and_open_door_to: room ID, rescue: 'Alice'	Proceed to provide your solution as a list of tuples in	
AS AS	KEY CONSTRAINTS:	chronological order.	
(Ĕ)	1. Each move must be between adjacent and connected rooms 2. Keys must be picked up before use 3.		
~	Locked doors require use of their specific key to unlock 4. Optimal path minimizes actions/distance 5. use_key	PROBLEM.	
	action always come right before unlock_and_open_door_to 6. If the response is missing any intermediate	EACTS:	
	action it is invalid - so it should include all the details necessary IMPORTANT: Use only provided IDs.	Room A6 and A5 are connected by an open door. Room A6 and	
	OUTPUT FORMAT REQUIREMENT:	B6 are connected by an open door. Room B6 and C6 are	
	Your solution must be formatted as a Python list of tuples representing each action in chronological order:	connected by an open door. Room C6 and D6 are connected by	
	[('start', 'RoomID'), ('move_to', 'RoomID'), ('pick_up_key', 'KeyID'),]	an open door. Room C5 and C4 are connected by an open door	
	Example format: [('start', 'A1'), ('move_to', 'B1'), ('pick_up_key', '3'), ('use_key', '3'),	Room C4 and D4 are connected by an open door. Room D6 and	
	('unlock_and_open_door_to', 'C1'), ('rescue', 'Alice')]	D5 are connected by a closed and locked door. The locked	Pro
	EXAMPLES:	door between D6 and D5 requires key 10. Key 10 is in room A5.	<u>p</u>
	INPLIT-	Room D6 and E6 are connected by an open door. Room D5 and	em
Examples	FACTS: Room C4 and C3 are connected by an open door. Room C3 and D3 are connected by an open door	D4 are connected by an open door. Room E6 and F6 are	T
	Room D5 and E5 are connected by an open door. Room A2 and A1 are connected by an open door. Room A3 and	connected by an open door. Room A4 and A3 are connected by	act
	B3 are connected by an open door. Room A1 and B1 are connected by an open door. Room A4 and A3 are	an open door. Bob is in room F6. Alice is in room C5.	S
	connected by an open door. Room E5 and E4 are connected by an open door. Room D4 and D3 are connected	YOUR SOLUTION:	
	by an open door. Room A5 and B5 are connected by an open door. Room D4 and E4 are connected by an open		
	door. Bob is in room D5. Alice is in room C4.		
	OUTPUT:		
	[('start', 'D5'), ('move_to', 'E5'), ('move_to', 'E4'), ('move_to', 'D4'), ('move_to', 'D3'), ('move_to', 'C3'),		
	('move_to', 'C4'), ('rescue', 'Alice')]		
	END OF EXAMPLES		

Figure 6: The complete prompt structure passed to the LLMs. This includes: Component 1 (System Instructions and Task Definition), one of the three Few-Shot Examples (Component 2, specifically a simple navigation task), Component 3 (Reasoning Guidance), and an illustration of where the Problem Instance Facts (Component 4) are inserted. For clarity and completeness, the full verbatim text for all three few-shot examples (Component 2) is provided in 7.

in Section 1 of the main paper and the discussion 1045 of failure modes in Section 2 of the main paper. 1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1059

1060

1061

1062

1063

1065

**Observed Violation Categories.** Failures in model solutions on seqBench tasks can be categorized into several types. Understanding these categories is crucial for interpreting model performance and failure modes. Key types of violations observed include:

- Adjacency errors (e.g., attempting to move between unconnected rooms).
- Locked door errors (e.g., navigating through locked doors without the correct key or without unlocking them).
- Key usage errors (e.g., attempting to use keys not yet collected, or using the wrong key for a door).
- · Path inefficiency (e.g., taking unnecessary detours or redundant actions; while not always a hard violation that stops progress, this contributes to solutions not matching the optimal path and thus failing Pass@1).

• Missed critical actions (e.g., failing to pick up 1066 a necessary key or unlock a required door). 1067 This is a key failure mode discussed in the 1068 main paper (Section 2.4) and is often reflected 1069 in metrics like low recall or a low progress 1070 ratio if the omission occurs early and prevents further correct steps. 1072

1071

1073

1074

1076

1077

1078

1079

Identifying these distinct categories of errors provides a more granular understanding of why models fail on sequential reasoning tasks and helps in the interpretation of aggregate performance metrics reported in the main paper.

#### **Violation Map: Qualitative Examples of B.3 Model Failures**

This section provides qualitative examples of char-1080 acteristic model failures to illustrate common error 1081 types. These examples visually support the dis-1082 cussion of failure modes in the main paper (Sec-1083 tion 2.4, "A Key Failure Mode: Omission of Criti-1084 cal Steps"). Figure 8 illustrates a significant error 1085 by Gemini-2.5-Flash on a complex task, where the 1086 model generates an illegal path, bypassing neces-1087 1. **Example 1 (Simple Navigation):** This example, as shown in Figure 6, involves navigating a maze with only open doors.

```
EXAMPLE:
INPUT:
Maze Structure: Room C4 and C3 are connected by an open door. Room C3 and D3 are
      connected by an open door. Room D5 and E5 are connected by an open door.
      Room A2 and A1 are connected by an open door. Room A3 and B3 are connected
      by an open door. Room A1 and B1 are connected by an open door. Room A4 and
      A3 are connected by an open door. Room E5 and E4 are connected by an open
      door. Room D4 and D3 are connected by an open door. Room A5 and B5 are
      connected by an open door. Room D4 and E4 are connected by an open door. Bob
      is in room D5. Alice is in room C4.
OUTPUT:
Solution: [('start', 'D5'), ('move_to', 'E5'), ('move_to', 'E4'), ('move_to', '
      D4'), ('move_to', 'D3'), ('move_to', 'C3'), ('move_to', 'C4'), ('rescue', '
      Alice')]
```

2. Example 2 (Single-Key Backtracking): This example introduces a single locked door and a corresponding key.

```
EXAMPLE:
INPUT:
Maze Structure: Room A1 and A2 are connected by an open door. Room A2 and B2 are
connected by an open door. Room B1 and B2 are connected by an open door.
Room B1 and C1 are connected by an open door. Room C1 and C2 are connected
by a closed and locked door. Door between C1 and C2 requires key 1. Key 1 is
in room A2. Bob is in room A1. Alice is in room C2.
OUTPUT:
Solution: [('start', 'A1'), ('move_to', 'A2'), ('pick_up_key', '1'), ('move_to',
'B2'), ('move_to', 'B1'), ('move_to', 'C1'), ('use_key', '1'), ('
unlock_and_open_door_to', 'C2'), ('move_to', 'C2'), ('rescue', 'Alice')]
```

3. Example 3 (Multi-Key Backtracking): This example presents a more complex scenario with multiple locked doors and keys, requiring more extensive backtracking.

```
EXAMPLE:
INPUT:
Maze Structure: Room B5 and B4 are connected by a closed and locked door. The
locked door between B5 and B4 requires key 3. Key 3 is in room B5. Room B5
and C5 are connected by a closed and locked door. The locked door between B5
and C5 requires key 16. Key 16 is in room C5. Room B4 and C4 are connected
by an open door. Room C4 and C3 are connected by an open door. Room C3 and
D3 are connected by a closed and locked door. The locked door between C3 and
D3 requires key 10. Key 10 is in room C4. Room D5 and D4 are connected by
an open door. Room D4 and D3 are connected by an open door. Room A5 and B5
are connected by an open door. Bob is in room C5. Alice is in room D5.
OUTPUT:
Solution: [('start', 'C5'), ('pick_up_key', '16'), ('use_key', '16'), ('
unlock_and_open_door_to', 'B5'), ('move_to', 'B5'), ('pick_up_key', '3'), ('
use_key', '3'), ('unlock_and_open_door_to', 'B4'), ('move_to', 'B4'), ('
move_to', 'C4'), ('pick_up_key', '10'), ('move_to', 'C3'), ('use_key', '10')
, ('unlock_and_open_door_to', 'D3'), ('move_to', 'D3'), ('move_to', 'D4'),
('move_to', 'D5'), ('rescue', 'Alice')]
```

Figure 7: Few-shot examples provided to guide the LLMs in the maze-solving task. These examples demonstrate simple navigation, single-key backtracking, and multi-key backtracking scenarios. The three examples illustrate increasing levels of complexity.

sary steps and locked doors. This exemplifies a breakdown in multi-step planning. Additionally,

1089

Figure 9 shows another common 'adjacency error,' where a model attempts to jump between uncon-

1092nected rooms. This type of error reveals a critical1093lapse in grounding its generated actions within the1094spatial adjacencies explicitly stated by the task's1095input facts.

#### B.4 Quantitative Analysis of Error Patterns

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

To understand how and when models begin to fail within a reasoning sequence, we analyze the distribution of the first violation step. We record the time step at which the initial violation occurs in a model's generated path. Aggregating this stepindexed data across multiple instances allows us to create temporal distributions of errors. These distributions help determine whether errors tend to cluster early in the reasoning process (potentially indicating issues with initial planning or understanding of the overall problem complexity) or accumulate later (suggesting difficulties in maintaining long chains of inference or context). This analysis complements the discussion in the main paper (Section 2.4, "Path-Length Dependent First Errors: The Burden of Anticipated Complexity").

Figure 10 shows how the distribution of these first-error positions shifts with the overall problem complexity, represented by logical depth (L). As detailed in the main paper, an increase in L tends to cause errors to occur earlier in the reasoning chain.

Similarly, Figure 11 illustrates how the introduction of contextual noise (distracting facts) affects the point of failure. Increased noise also tends to precipitate earlier errors in the reasoning sequence, as discussed in the main paper in relation to sensitivity to noise (Section 2.3) and its impact on error patterns (Section 2.4).

#### C Supplementary Figures

This appendix provides supplementary figures that offer further visual support for analyses presented in the main paper. These figures illustrate the impact of various complexity dimensions and provide comparative views of model performance, elaborating on points made throughout Section 2 (Benchmarking Results) of the main paper.

Figure 12 details the performance of Llama-4 Maverick-17B-128E-Instruct under varying levels of noise and fact shuffling. This supports the discussion in the main paper (Section 2.3, on how these factors, especially in combination, affect success rates, with noise being a dominant factor.

To illustrate the performance consistency and disparities across different models, as detailed in

Section 2.6, Figure 13 presents scatter and density 1141 plots of mean progress ratios. These plots clearly 1142 demonstrate that model performance hierarchies 1143 are not strictly linear. They reveal 'performance 1144 inversions'-instances, also noted in Section 2.6, 1145 where models with typically lower overall perfor-1146 mance (e.g., lower average  $L_0$ ) occasionally solve 1147 specific complex problems that models with higher 1148 average  $L_0$  values fail on. 1149

1150

1151

1152

1153

1154

1155

1156

1157

Figure 14 isolates the impact of shuffle ratio on model performance when other factors like noise are controlled. This visualization corresponds to the findings discussed in the main paper (Section 2.3, "Fact Ordering (Shuffle Ratio)") that simple reordering of facts has a minimal impact on the performance of the evaluated models under lownoise conditions.

#### **Optimal Path**

#### Model Path





Figure 8: Illustrative failure case for Gemini-2.5-Flash on a 40x40 task. Left: Optimal path (yellow). Right: Model's generated path showing an illegal adjacency jump (red arrow), bypassing multiple rooms and a locked door, despite only supporting facts being provided. This highlights a breakdown in multi-step planning.



Figure 9: Illustrative failure case of an 'adjacency error' in model-generated pathfinding. The left panel displays the optimal path (yellow) to the target (triangle). The right panel shows a suboptimal path (purple) generated by the model. This example highlights a common error where, after a sequence of actions (in this scenario, following a key acquisition), the model fails to navigate through valid connections. Instead, it attempts to 'jump' directly between two unconnected rooms. This violation of room adjacency constraints is a key challenge in model performance.

Solution steps: 20

Solution steps: 60 Solution steps: 100 Solution steps: 140 Solution steps: 180 Solution steps: 220 Solution steps: 260 Solution steps: 300 0 50 100 150 200 250 300 max progress step

Figure 10: Distribution of first-violation steps for Gemini-2.5-Flash across varying logical depths (L). As L (total required path length) increases, the distribution of first errors tends to shift leftward, indicating that models are more likely to fail at earlier steps in longer problems. This suggests that anticipated global complexity impacts reasoning from the outset.



Figure 11: Impact of increasing noise ratio on the distribution of failure steps for Gemini 2.5 Flash. As noise (proportion of distracting facts) increases, failures tend to occur earlier in the reasoning chain. This reflects increased difficulty in isolating relevant information and maintaining focus.



Figure 12: Pass@1 success rate for Llama-4 Maverick-17B-128E-Instruct versus solution length (L) under different noise and shuffle ratios. Left: Linear scale. Right: Log-linear scale. Performance degrades with increased noise but is less affected by shuffle ratios.



Figure 13: Scatter and density plots of progress ratios per task instance, comparing model pairs on the tasks. These plots illustrate performance agreement and disparities on the same instances of pathfinding tasks. Notably, Gemini-2.5-Flash (example) often succeeds on instances where other models achieve near-zero progress. Data from experiments in Figure 1 (main paper).



Figure 14: Impact of shuffle ratio on Pass@1 success rate. Varying the degree of mixing (shuffle) between supporting and distracting facts shows minimal impact on performance for Gemini 2.5 Flash and Llama-4 Maverick, suggesting robustness to fact order when noise is controlled. The generation and sampling of maze instances for these tasks follow the same methodology detailed for experiments in the main paper (Figures 3 and 4).