WirelessMathBench: A Mathematical Modeling Benchmark for LLMs in Wireless Communications

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have achieved impressive results across a broad array of tasks, yet their capacity for complex, domain-specific mathematical reasoning-particularly in wireless communications-remains underexplored. In this work, we introduce WirelessMathBench a novel benchmark specifically designed to evaluate LLMs on mathematical modeling challenges to wireless communications engineering. Our benchmark consists of 587 meticulously curated questions sourced from 40 state-of-theart research papers, encompassing a diverse spectrum of tasks ranging from basic multiplechoice questions to complex equation completion tasks, including both partial and full completions, all of which rigorously adhere to phys-017 ical and dimensional constraints. Through extensive experimentation with leading LLMs, we observe that while many models excel in basic recall tasks, their performance degrades significantly when reconstructing partially or fully obscured equations, exposing fundamental limitations in current LLMs. Even DeepSeek-R1, the best performer on our benchmark, achieves an average accuracy of only 38.05%, with a mere 7.83% success rate in full equation com-027 pletion. By publicly releasing WirelessMath-Bench along with the evaluation toolkit, we aim to advance the development of more robust, domain-aware LLMs for wireless system analysis and broader engineering applications.

1 Introduction

034

042

Large Language Models (LLMs) have recently demonstrated groundbreaking performance across a diverse range of natural language tasks—from general language understanding (Brown et al., 2020; Wang et al., 2019b,a) and code generation (Lu et al., 2021) to elementary mathematical reasoning (Cobbe et al., 2021; Hendrycks et al., 2021a). Advanced models such as OpenAIo1 (OpenAI, 2024) and DeepSeek-R1 (Guo et al.,



Figure 1: Example task from WirelessMathBench a system model derivation from wireless communications literature. The derivation progresses from a multiple-choice question to progressive mask completion questions, and finally to the full formula derivation, testing the model's ability to reason through complex channel reflections and matrix operations.

2025) have further extended these capabilities, especially when supplemented with chain-of-thought strategies that enable clear, step-by-step solution processes. Nevertheless, despite these notable achievements, current state-of-the-art LLMs still encounter significant difficulties when tackling highly intricate problem statements. In particular, tasks that demand deep conceptual insights, rigorous validation of physical feasibility, and the careful management of tightly interrelated parameter sets continue to pose formidable challenges

053

043

101

102

103

105

(Mirzadeh et al., 2025; Zhang et al., 2023; He et al., 2024).

In many engineering fields-wireless communications in particular-mathematical modeling is indispensable. The design and analysis of modern wireless systems require not only accurate numerical computation but also precise symbolic derivations that honor strict physical and dimensional constraints. Tasks such as channel estimation (Yin et al., 2013; Liu et al., 2022; An et al., 2023), beamforming (Chu et al., 2023; Spencer et al., 2004), and multi-antenna system design (Huang et al., 2020; An et al., 2024; Zheng et al., 2024) involve intricate matrix operations, multi-stage derivations, and domain-specific lexicon. Even minor errors in symbolic manipulation can lead to significant performance degradation or non-compliance with industry standards (Bjornson et al., 2013).

Although recent work has leveraged LLMs for technical definition retrieval in wireless communications (Shao et al., 2024; Maatouk et al., 2023; Zou et al., 2024; Maatouk et al., 2024), few studies have directly addressed the challenges associated with multi-step derivations and symbolic manipulation in this specialized domain. This observation raises a broader question: *To what extent are LLMs capable of emulating the mathematical derivations and analytical typically by an engineer or researcher in the field of real wireless communications*?

To bridge this gap, we introduce WirelessMath-Bench a comprehensive benchmark specifically designed to test LLMs on the real-work wireless engineering orientation mathematical reasoning. WirelessMathBench comprises 587 high-quality questions sourced from 40 state-of-the-art papers, each carefully annotated and validated by domain experts to ensure accuracy. These questions span a variety of system models (e.g., Multiple-Input and Multiple-Output (MIMO), Non-orthogonal multiple access (NOMA), Reconfigurable Intelligent Surfaces (RIS)) and problem settings (e.g., channel estimation, beamforming), encompassing multiplechoice, fill-in-the-blank, and open-ended questions at various levels. Table1 highlights key differences between WirelessMathBench and other math benchmarks, ours is the only dataset of expert difficulty level and contains real-world engineering problems. Figure 1 illustrates how a single math formula escalates from a basic multiple-choice query to a fully masked equation derivation, reflecting the complexity of real-world wireless system

analysis.

Our extensive experiments show that while leading LLMs perform well on simpler tasks (e.g., multiple-choice questions with over 75% accuracy), their performance drops dramatically on advanced derivation tasks (progressive masking and complete equations). Even the strongest model we evaluated, DeepSeek-R1 (Guo et al., 2025), only manages a 7.83% success rate in fully masked derivations, underscoring a fundamental gap between current LLM capabilities and the complex demands of wireless systems analysis.

By publicly releasing WirelessMathBench along with its evaluation toolkit, we aim to spur progress toward LLMs that are not only fluent in natural language but also capable of rigorous, domainspecific mathematical reasoning. We envision that WirelessMathBench will serve as a catalyst for innovation in mathematical reasoning capabilities, domain-adaptive pre-training techniques, and advanced thought-chaining strategies, ultimately propelling LLMs toward more robust scientific and engineering problem-solving.

2 Related Work

General-Purpose LLM Benchmarks. In recent years, rapid advancements in LLMs—exemplified by models such as GPT-3 (Brown et al., 2020), GPT-4 (Achiam et al., 2023), LLaMA (Touvron et al., 2023), Gemini (Team et al., 2023), and DeepSeek-R1 (Guo et al., 2025)—have spurred extensive evaluations on benchmarks like GLUE (Wang et al., 2019b), SuperGLUE (Wang et al., 2019a), and GSM8K (Cobbe et al., 2021). However, despite covering a broad spectrum of linguistic tasks, they typically lack the depth and specificity required to evaluate rigorous mathematical modeling or the domain-specific symbolic reasoning needed for complex technical applications.

Mathematical Reasoning Benchmarks. A parallel research stream has focused on the mathematical and symbolic reasoning abilities of LLMs. Early mathematical benchmarks (Amini et al., 2019; Cobbe et al., 2021; Koncel-Kedziorski et al., 2016; Ling et al., 2017; Hendrycks et al., 2021b) evaluate models on elementary arithmetic, algebra, and calculus problems. Recently, as the complexity of the problem increases, some benchmarks introduce competitionlevel problems that combine mathematical logic and background knowledge (Yu et al., 2024; 106 107 108

109

110

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

154

Benchmark	Diffuculty Level	QuestionType	Real Engineering Tasks	#Test Size
GSM8K (Cobbe et al., 2021)	Elementary School	OE	No	1,319
MATH (Hendrycks et al., 2021b)	High School	OE	No	5,000
OCWCourses (Lewkowycz et al., 2022)	University	OE	No	272
MMMU (Yue et al., 2024)	University	MC,OE	No	1983
OlympiadBench (He et al., 2024)	Competition	OE	No	8,476
WirelessMathBench	Expert	MC, FB, OE	Yes	587

Table 1: Comparison of representative mathematical benchmarks with **WirelessMathBench**. Existing datasets largely focus on elementary, high school, or Olympiad-level problems in purely theoretical contexts, while **Wire-lessMathBench** targets real-world, expert-level engineering tasks under strict dimensional and physical constraints. We note that open-ended (OE) tasks typically require free-form answers, MC indicates multiple-choice, and FB refers to fill-in-the-blank.

Hendrycks et al., 2021a; Arora et al., 2023; Frieder et al., 2024). For more advanced mathematical reasoning, datasets like MMMU (Yue et al., 2024), OCWCourses (Lewkowycz et al., 2022) and U-MATH (Chernyshev et al., 2024) focuses on university-level mathematics problems. MiniF2F (Zheng et al., 2022), AlphaGeometry (Trinh et al., 2024), OlympiadBench (He et al., 2024), and MathOdyssey (Fang et al., 2024) go further to Olympiad-level problems that require more advanced mathematical reasoning. Yet, these datasets do not capture the unique constraints or specialized notations found in applied domains like wireless communications.

156

157

158

159

160

161

163

164

165

166

167

168

169

Domain-Specific Benchmarks. To overcome the 170 limitations of general-purpose evaluations, sev-171 eral domain-specific benchmarks have been devel-172 oped for tasks that demand technical precision and 173 specialized reasoning. For example, customized 174 175 benchmarks have been developed for legal document analysis (Guha et al., 2024), chemical prop-176 erty inference (Guo et al., 2023), and scientific 177 reasoning (Lu et al., 2022; Wang et al., 2024; Sun 178 et al., 2024). To evaluate LLMs in more specialized 179 domain tasks, recent works have introduced bench-180 marks like MLAgentBenchmark (Huang et al., 181 2024), which evaluates LLMs' ability to solve ma-182 chine learning tasks, AI-Researcher (Si et al., 2025) evaluate can LLMs generate research ideas, and 184 SWE-Bench (Jimenez et al., 2024) evaluate LLMs' 185 ability to solve real-world software engineering 186 tasks. These studies highlight that an in-depth evaluation of LLMs in specialized fields reveals 188 that LLMs have strong potential in different professional fields.

191LLMs in Wireless Communications. Wireless192communications impose stringent requirements on193mathematical precision, particularly for tasks such

as channel estimation, interference management, and beamforming (Cadambe and Jafar, 2008; Shi et al., 2011; Gesbert et al., 2010). Some preliminary works have explored the use of LLMs in wireless contexts, focusing on domain-specific knowledge extraction and basic recall of technical standards (Maatouk et al., 2023; Shao et al., 2024; Maatouk et al., 2024). Notably, TelecomGPT (Zou et al., 2024) has extended LLM capabilities to higher-level tasks like wireless-specific code generation and formula completion. However, these early works usually emphasize knowledge retrieval or summarization, without considering testing what tasks LLMs can accomplish in actual wireless communication engineering systems. 194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

229

230

231

In this work, we introduce WirelessMathBench to address these gaps. Unlike existing wireless or purely mathematical benchmarks, Wireless-MathBench offers tasks that systematically combine multiple-choice questions with progressively masked formula derivations, all drawn from stateof-the-art research papers. The goal is to evaluate both symbolic reasoning and domain knowledge under realistic conditions, capturing the nuanced interplay of mathematical derivations and physical feasibility inherent in wireless communications. By providing a diverse set of tasks and domaininformed evaluation metrics, WirelessMathBench aims to facilitate collaborative advances in both LLMs and wireless communication engineering, ultimately enabling more powerful AI-assisted solutions for next-generation wireless networks.

3 The WirelessMathBench Benchmark

In this section, we present WirelessMathBench, a new benchmark specifically designed to evaluate LLMs on mathematical modeling tasks within wireless communications. We begin by discussing the rationale behind our benchmark design (Sec-



Figure 2: Overview of the data collection and annotation pipeline for WirelessMathBenchThe process involves selecting high-quality research papers, extracting system models from papers, curating tasks of varying complexity levels, and reviewing each task for clarity and correctness.

tion 3.1), followed by the details of our data collection and annotation pipeline (Section 3.2). We then explain how we construct questions of varying complexity levels, as well as our progressive masking methodology (Section 3.3).

3.1 Design Principles

232

236

239

240

241

243

245

247

248

251

265

The creation of WirelessMathBench is motivated by two core observations. First, recent work shows that LLMs can effectively assist humans in highly specialized tasks (Guha et al., 2024; Guo et al., 2023; Lu et al., 2022), underscoring their potential when provided with sufficient domain context. Second, LLMs have demonstrated the capacity to handle increasingly difficult mathematics, including Olympiad-level challenges (He et al., 2024; Fang et al., 2024). These findings suggest a substantial opportunity to push the limits of LLMs in areas where complex, domain-specific mathematics—such as wireless communications—plays a central role.

Building on these insights, WirelessMathBench is designed around two key principles:

- 1. **Real-World Complexity.** Each task is sourced directly from peer-reviewed research, reflecting the authentic modeling challenges faced in wireless systems.
- 2. **Multi-Tiered Progression.** Tasks range from basic multiple-choice questions to fully masked derivations, providing graduated levels of difficulty that capture both foundational knowledge and advanced reasoning.

3.2 Data Collection and Annotation

As illustrated in Figure 2, the data collection and annotation process for WirelessMathBench involves four main steps: paper selection, system model extraction, task curation, and domain expert review. 266

267

268

269

270

271

272

273

274

275

276

277

278

281

284

285

288

289

290

291

292

294

295

297

298

299

300

301

302

Paper Selection and Coverage. To capture the authentic complexities of wireless communications, we begin by identifying high-impact papers from top-tier publication venues that are freely accessible on arXiv. Table 2 summarizes the coverage of model-based and problem-based categories; in total, we select 40 papers spanning core techniques (e.g., MIMO, NOMA, RIS) and research focuses (e.g., channel estimation, beamforming). We emphasize works that feature nontrivial mathematical derivations—such as optimization formulations and multi-stage channel modeling-over those limited to empirical or simulation-based heuristics. A summary of the high-frequency keywords across the dataset is shown in Figure 3, highlighting the diverse wireless communication topics covered in the benchmark. We aim to include tasks that reflect the symbolic depth and physical constraints that are indispensable for real-world wireless engineering and state-of-the-art wireless research.

System Model Extraction. Our pipeline starts by applying a specially designed LLM template that systematically scans each research paper, isolating key mathematical expressions and relevant contextual descriptions. This initial extraction step is semi-automated: LLMs produce a structured draft containing a concise overview of the system model, assumptions, and principal formulas. Subsequently, domain experts review and refine these drafts, ensuring that the extracted material is both accurate (i.e., symbolically consistent with the original text) and self-contained (i.e., providing enough background to be understood independently). This hybrid method combines the scalability of automated extraction with meticulous expert checks,

Category Type	Topic Category	Number of Papers
	RIS	19
	MIMO	12
	UAV	6
Model-based	ISAC	6
	Satellite	4
	SIM	3
	NOMA	2
	Beamforming	18
	Channel Estimation	12
Droblom based	Performance Analysis	8
Problem-based	Trajectory Design	5
	Power Allocation	5
	Resource Management	4
Total		40

Table 2: Distribution of the WirelessMathBench benchmark papers according to model-based and problembased categories, along with their respective topic areas. A total of 40 papers are included, covering key themes in wireless communications. Note that some papers may span multiple topic categories.

ensuring the resulting text is accurate, symbolically consistent, and sufficiently self-contained for sub-sequent tasks.

303

304

305

306

310

311

312

313

314

315

319

321

323

324

326

327

328

Task Curation. Once the system models have been extracted, we systematically transform them into precise question-answer pairs via a three-step process:

- 1. **Identify Core Equations:** We select representative formulas from each paper—such as channel gain expressions, optimization objectives, or multi-hop path-loss derivations—that encapsulate critical wireless engineering challenges.
- Construct Questions: For each core equation, we generate questions at varying levels of difficulty. This includes: (a) multiple-choice questions targeting definitions or partial operations;
 (b) progressively masked fill-in-the-blank questions that require incremental reasoning; and
 (c) full equation completion tasks that demand derivation of the entire expression.
- 3. Annotate and Review: Each question is accompanied by contextual notes, and domain experts validate correctness and clarity. When necessary, we refine notation or provide short explanations to ensure the questions are self-contained and can be tackled without external references

329Domain Expert Review.Lastly, all questions330undergo a multi-round review by senior wireless331researchers.They verify notation accuracy and



Figure 3: A word cloud illustrating the most frequent keywords in the WirelessMathBench benchmark, which reflects the range of wireless communication topics covered.

domain applicability (e.g., check for appropriate dimensionality, and coherent modeling assumptions), and remove any ambiguous or misleading content. The remaining problems constitute the final WirelessMathBench dataset: a set of carefully selected tasks that embody typical mathematical derivations in advanced wireless communication engineering and research.

3.3 Task Design and Masking Strategies

To full evaluate the capabilities of LLMs at different levels of difficulty in mathematical modeling of wireless communications, WirelessMathBench incorporates three distinct task types. Each question leverages real-world system equations derived from state-of-the-art research papers, ensuring that the benchmark reflects both conceptual diversity and practical engineering relevance. At the same time, each independent question is accompanied by a brief description of the relevant wireless scenario (e.g., UAV relay or multi-antenna beamforming), providing the necessary domain and scenario background information.

Multiple-Choice Questions (MCQs). These questions require the solver to select the correct mathematical expression from a set of closely related distractors, with each MCQ carefully designed to test the model's ability to recognize and recall key elements of wireless system modeling. For example, a typical MCQ may present several equations for a wireless channel, of which only one formula satisfies both the correct dimensions and the physical constraints of the system under consideration.

Progressively masked fill-in-the-blank questions. In this task, a system model formula is progressively presented in a partially masked form across

367

332

333

334

335

Model	Source	Size
OpenAI-o1	(OpenAI, 2024)	unknown
DeepSeek-R1	(Guo et al., 2025)	671B
GPT-Family	(OpenAI, 2024; Achiam et al., 2023)	unknown
DeepSeek-V3	(DeepSeek-AI et al., 2024)	671B
Gemini-Famliy	(Google DeepMind, 2024)	unknown
Qwen2.5-Math	(Yang et al., 2024)	7B, 72B
LLaMA	(Grattafiori et al., 2024)	8B, 70B
LLaMA-3-8B-Tele	(Maatouk et al., 2024)	8B
Mistral-7B	(Jiang et al., 2023)	7B

Table 3: LLMs evaluated on WirelessMathBench.

368three different masking levels. Each progressively369masked instance is treated as an independent sub-370problem, requiring the model to infer and recon-371struct the missing information at each stage. The372masking levels range from isolated single-variable373omissions to multi-variable occlusions, with vary-374ing degrees of accompanying prompt text to pro-375vide contextual guidance.

Full Equation Completion (FEC) question. For 376 the most challenging question, the full equation is entirely hidden. The solver is provided with only a succinct description of the wireless scenario 379 (for example, a base station-relay-user link with specific path loss characteristics) and must derive the complete expression from first principles. This task assesses the model's ability to reconstruct the entire derivation-from fundamental definitions 384 (like channel gain or fading coefficients) to the final expression-while ensuring dimensional accuracy and adherence to domain-specific constraints (such as path loss exponent and transmit power limits). It 388 represents the level of performance expected from a human expert in wireless communications. 390

> In summary, by combining MCQs, progressively masked tasks, and full equation completions, WirelessMathBench offers a comprehensive, finegrained evaluation of a model's capability to perform both symbolic reasoning and domain-specific derivations in wireless communications.

4 Experiments

391

395

398

400

401 402

403

404

405

406

We now present a comprehensive evaluation of WirelessMathBench, focusing on how leading LLMs handle wireless communications–specific mathematical modeling. We begin by detailing our experimental setup (Section 4.1), including prompt design and model baselines, then discuss our main results (Section 4.2), and conclude with an error analysis that highlights key challenges (Section 4.3).

4.1 Experiment Setup

Evaluation Workflow. All experiments are conducted in a zero-shot setting using unified prompt templates across different question types for consistent evaluation. For each task in WirelessMath-Bench, we provide the corresponding prompt to each model and collect the answers it generates. Our evaluation pipeline is now completed in two main ways. For multiple-choice questions, we directly extract the output answers and compare their consistency with the annotation results. For Progressive Masking Filling and Fully maded questions, since polynomials may have a certain number of possible answers, we use the help of LLMs (GPT-40 is selected in our experiment) to complete the evaluation, similar to(Fang et al., 2024; Chernyshev et al., 2024). The overall performance is reported as the average accuracy of all tasks. Detailed prompt examples and scoring criteria are provided in the Appendix B.

Baselines. Table 3 (in the main text) lists the principal models tested. We include leading reasoning models(e.g., DeepSeek-R1, OpenAI-o1), large-scale general-purpose LLMs (e.g., GPT-4, Gemini), and specialized models (e.g., Qwen2.5-Math) to capture a broad range of capabilities. For open-source models like LLaMA, we also explore domain-specific variants trained on a telecom corpus (e.g., LLaMA-3-8B-Tele) to gauge the benefit of targeted adaptation. All hyperparameters follow each model's respective default or recommended settings, and no additional chain-of-thought prompting is provided beyond the standard instructions above.

4.2 Main Results

Table 4 presents the performance of sixteen LLMs across five metrics in WirelessMathBench: (1) Multiple-choice Question (MCQ) accuracy, (2–4) progressive masking fill-in at three difficulty levels (Level 1, Level 2, Level 3), (5) Full Equation Completions (FEC), and the overall average accuracy (Avg. Acc). Our key findings are summarized below:

Reasoning-Oriented Models Show Advantages.450Models that incorporate explicit chain-of-thought451or advanced reasoning techniques—like DeepSeek-452R1 and OpenAI-o1—consistently outperform sim-453pler large-scale baselines. The average accuracy of454DeepSeek-R1 is 38.05%, and the average accuracy455

443

444

445

446

447

448

449

407

408

409

410

411

412

413

Model	MCO	Progressive Masking Filling		ressive Masking Filling FEC Avg. Acc		
	nic y	Level 1	Level 2	Level 3	120	11,8,1100
DeepSeek-R1	76.00%	60.00%	34.91%	12.50%	7.83%	38.05%
OpenAI-o1	66.40%	59.17%	32.17%	8.04%	6.96%	34.55%
OpenAI-o1-mini	66.40%	53.33%	29.57%	10.71%	4.35%	32.87%
GPT-40	72.80%	42.50%	28.70%	6.25%	4.35%	30.92%
GPT-4	53.60%	38.33%	18.26%	3.57%	4.35%	23.62%
GPT-3.5-turbo	45.60%	7.50%	10.43%	1.79%	1.74%	13.41%
DeepSeek-V3	78.40%	50.00%	24.35%	6.25%	6.96%	33.19%
Gemini-2.0-flash	71.20%	40.83%	24.35%	5.36%	4.35%	29.22%
Gemini-1.5-pro	65.60%	43.33%	29.57%	9.82%	6.09%	30.88%
Gemini-1.5-flash	66.40%	37.50%	13.91%	2.68%	4.35%	24.97%
Qwen2.5-Math-72B	70.40%	37.50%	26.09%	7.14%	6.09%	29.44%
LLaMA-3.3-70B	65.60%	38.33%	17.39%	2.68%	6.09%	26.02%
Qwen2.5-Math-7B	58.40%	21.67%	6.96%	4.46%	1.74%	18.82%
LLaMA-3-8B-Tele	40.80%	11.67%	4.35%	2.68%	0.87%	12.07%
LLaMA-3-8B	45.60%	10.83%	7.83%	2.68%	2.61%	13.91%
Mistral 7B	38.40%	20.00%	4.35%	0.89%	0.87%	12.90%

Table 4: Experimental results of state-of-the-art LLMs on WirelessMathBench. The table shows the performance of each model on MCQ, progressively masked filling and full equation completion tasks.

of **OpenAI-o1** is 34.55%, while the accuracy of other large-parameter models hovers around 30%.
This suggests that explicit reasoning strategies contribute substantially to managing multi-step symbolic derivations in wireless communications tasks.

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

Strong MCQ Performance but Rapid Decline in Derivations. Several models, including DeepSeek R1, V3, GPT-4, and Gemini-2.0, exceed 70% accuracy on MCQs, showing that they can find the correct formula given background knowledge and given error options, indicating that they can understand the modeling process and matrix operations in the communication domain to some extent. However, these MCQ gains generally do not extend to more complex derivation tasks, where most models' accuracy falls dramatically. For instance, **DeepSeek-V3** achieves the highest MCQ score at 78.40%, drop to 6.25% in Level 3 masking filling, and 6.96% in FEC.

475 Progressive Masking Emphasizes Multi-Step
476 Reasoning Gaps. When forced to reconstruct
477 partially hidden expressions, model performance
478 declines in proportion to the level of masking.
479 When forced to reconstruct partially hidden expressions, model performance degrades with in480 pressions, model performance degrades with in481 creasing levels of masking. Models with implicit

reasoning logic significantly outperform the others, with **DeepSeek-R1** in particular leading on these tasks—achieving 60.00% at Level 1 and 33.91% at Level 2, suggesting more robust chaining of thoughts. However, even DeepSeek-R1 struggles at level 3 (12.50%), highlighting the difficulty of maintaining symbolic coherence under heavily ambiguous conditions.

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

Fully Masked Equation Completion Remains Challenging. Most models attain only singledigit accuracy (2–7%) in the FEC task, where the entire equation is hidden. **DeepSeek-R1**'s 7.83% and **OpenAI-o1**'s 6.96% are the best in this category, but both remain low in absolute terms, indicating that fully reconstructing multi-step derivations without partial clues poses a significant challenge.

Domain-Focused Models Show Improvements. Models that are specifically tuned for mathematical reasoning—such as Qwen2.5-Math—demonstrate improved performance over other models with a similar parameter count, both in terms of overall average accuracy and on individual subtasks. In particular, **Qwen2.5-Math-72B** achieves an average accuracy of 29.44%, which is on par with the average performance of most commercial models. However, fine-tuning general-purpose models like



Figure 4: Error distribution among 40 annotated DeepSeek-R1 errors.

508 LLaMA to telecom-specific data (e.g., LLaMA-3-8B-Tele) yields only limited benefits. This is likely because the telecom fine-tuning data predominantly 510 consists of wireless protocols, whereas the prob-511 lems in WirelessMathBench require handling long 512 contexts and performing high-level mathematical 513 514 reasoning.

4.3 Error Analysis

515

524

529

To better understand the limitations of the model-516 generated solutions, we randomly sampled and re-517 viewed 40 failure answers by DeepSeek-R1 on the progressively masked filling and fully masked equa-519 tion tasks, grouping them into several recurring categories. Figure 4 summarizes the distribution of 521 these errors.

Partial Fill Mismatch (31%). A prevalent challenge in progressive masking tasks is that the model often merges multiple placeholders into a single bracket or places the correct term in the 526 wrong mask position. For instance, it may combine $\sqrt{pK}\beta_{mk}\mathbf{y}_{pm}^{H}\boldsymbol{\phi}_{k}$ into one placeholder even 528 though the prompt explicitly requests separate terms $(\sqrt{pK}\beta_{mk} \text{ and } \mathbf{y}_{pm}^{H}\boldsymbol{\phi}_{k})$ across two different 530 masks.

Symbol Misinterpretation (29%). This type of error arises when the model chooses the wrong 533 symbol or omits key symbolic elements in the final equation. An example is substituting \mathbf{H}_{BR} instead 535 of $\mathbf{H}_{\mathrm{BR}}^{H}$ in a channel derivation.

Incorrect Equation Derivation (24%). Several 538 solutions fail to follow the correct derivation path, either missing crucial intermediate steps or injecting extraneous components. In longer sequences, a single early mistake (e.g., confusing pilot power 541 p with user transmit power ρ_k) tends to propa-542

gate, causing the final expression to be structurally flawed despite appearing superficially similar.

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

Irrelevant System Mixing (11%). We also observe instances in which the model introduces extraneous terms or assumes mismatched system settings. For example, it might inject NOMA-based interference factors into an RIS-MIMO scenario with no multi-user superposition, or switch to an entirely different beamforming constraint unconnected to the original problem statement.

Other Errors (4%). A small fraction of errors are too context-specific to classify under the main categories. These include incomplete expressionswhere the answer stops abruptly without filling the entire formula-and mismatched expressions that repeat known placeholders without substituting the correct variables.

Overall, while the majority of failures fall into coherent error patterns, it is evident that the model struggles when tasks require nuanced symbol-tomask matching or integrative domain reasoning. Therefore, strengthening the model's ability to accurately derive reasoning and align domain knowledge is a key challenge for future improvements.

5 Conclusion

We introduced WirelessMathBench, the first benchmark that comprehensively evaluates LLMs' abilities to domain-specific mathematical engineering tasks in wireless communications. By presenting a broad range of tasks-from multiple-choice questions to progressively masked fill-in-the-blank and full equation completions-WirelessMathBench assesses how well models handle real-world wireless communications mathematical engineering challenges. Our experiments show that, although many leading LLMs perform well on simpler question types, their ability to reconstruct equations deteriorates significantly when partial or full derivation is required, exposing a critical shortfall in current LLM-assisted scientific innovation. Future work will expand the scope and complexity of these wireless challenges, with the aim of enhancing LLMs' mathematical reasoning and domain adaptation. By advancing their integration into the next-generation wireless systems, we ultimately strive toward the development of more capable, general-purpose AI solutions for scientific and engineering applications.

Ethical Considerations

Limitations

This paper focuses on the development of a bench-

mark for evaluating language models on mathemat-

ical modeling tasks in wireless communications.

The source data of WirelessMathBench is curated

from open-access research papers, ensuring that

the benchmark is built on publicly available information. Meanwhile, we resummarize the papers

and anonymize the content to prevent any poten-

tial privacy concerns. In experiments, we follow

all licensing agreements and terms of service for

the models evaluated, ensuring that our work is conducted in compliance with ethical guidelines.

While WirelessMathBench provides a comprehensive evaluation of LLMs on wireless mathematical

modeling tasks, several limitations remain. First,

it mainly covers text-based problems (e.g., sym-

bolic derivations), missing other key data types

like antenna diagrams, simulation plots, and Ra-

dio frequency (RF) measurements measurements,

which are crucial for real-world wireless tasks. Sec-

ond, while WirelessMathBench spans topics from

MIMO to RIS, it may not cover all emerging ar-

eas, such as quantum communication or terahertz

systems. Third, our automated evaluation checks

the final symbolic equivalence and dimensionality

plausibility but may miss incorrect reasoning at in-

termediate steps. Lastly, all tests were done in a

zero-shot setting. While this reflects real-world use,

it does not explore whether fine-tuning or retrieval-

based methods could improve results. Future ver-

sions of WirelessMathBench may include training

splits to support domain adaptation and wireless-

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.

Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Ha-

jishirzi. 2019. MathQA: Towards interpretable math

word problem solving with operation-based for-

malisms. In Proceedings of the 2019 Conference

of the North American Chapter of the Association for

Computational Linguistics: Human Language Tech-

nologies, Volume 1 (Long and Short Papers), pages

arXiv preprint arXiv:2303.08774.

specific fine-tuning.

References

- 604

- 607
- 608
- 611
- 612
- 613

614

- 615
- 617
- 618 619

624

- 631

632

633

636

637

2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

- Jiancheng An, Chau Yuen, Chongwen Huang, Mérouane Debbah, H. Vincent Poor, and Lajos Hanzo. 2023. A tutorial on holographic MIMO communications-part i: Channel modeling and channel estimation. IEEE Communications Letters, 27(7):1664-1668.
- Jiancheng An, Chau Yuen, Chao Xu, Hongbin Li, Derrick Wing Kwan Ng, Marco Di Renzo, Mérouane Debbah, and Lajos Hanzo. 2024. Stacked intelligent metasurface-aided MIMO transceiver design. IEEE Wireless Communications, 31(4):123–131.
- Daman Arora, Himanshu Gaurav Singh, and Mausam . 2023. Have LLMs advanced enough? a challenging problem solving benchmark for large language models. In The 2023 Conference on Empirical Methods in Natural Language Processing.
- Emil Bjornson, Michail Matthaiou, and Merouane Debbah. 2013. A new look at dual-hop relaying: Performance limits with hardware impairments. IEEE Transactions on Communications, 61(11):4512-4525.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems, 33:1877-1901.
- Viveck R. Cadambe and Syed Ali Jafar. 2008. Interference alignment and degrees of freedom of the K-user interference channel. IEEE Transactions on Information Theory, 54(8):3425-3441.
- Konstantin Chernyshev, Vitaliy Polshkov, Ekaterina Artemova, Alex Myasnikov, Vlad Stepanov, Alexei Miasnikov, and Sergei Tilga. 2024. U-MATH: A university-level benchmark for evaluating mathematical skills in llms. arXiv preprint arXiv:2412.03205.
- Jinjin Chu, Rang Liu, Ming Li, Yang Liu, and Qian Liu. 2023. Joint secure transmit beamforming designs for integrated sensing and communication systems. IEEE Transactions on Vehicular Technology, 72(4):4778-4791.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168.
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, and et al. 2024. Deepseekv3 technical report. Preprint, arXiv:2412.19437.

807

808

809

753

Meng Fang, Xiangpeng Wan, Fei Lu, Fei Xing, and Kai Zou. 2024. MathOdyssey: Benchmarking mathematical problem-solving skills in large language models using odyssey math data. *arXiv preprint arXiv:2406.18321*.

699

705

706

707

708

710

711

712

713

714

715

716

718

719

720

721

725

728

729

730

731

732

733

734

736

737

738

740

741 742

743

744

745

746

747

748

749

751

- Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Petersen, and Julius Berner. 2024. Mathematical capabilities of chatgpt. *Advances in neural information processing systems*, 36.
- David Gesbert, Stephen Hanly, Howard Huang, Shlomo Shamai Shitz, Osvaldo Simeone, and Wei Yu. 2010. Multi-cell MIMO cooperative networks: A new look at interference. *IEEE Journal on Selected Areas in Communications*, 28(9):1380–1408.
- Google DeepMind. 2024. Introducing gemini 2.0: our new ai model for the agentic era. Accessed: 2024-12-11.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, and et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré, Adam Chilton, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zambrano, et al. 2024. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. Advances in Neural Information Processing Systems, 36.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Taicheng Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh Chawla, Olaf Wiest, Xiangliang Zhang, et al. 2023. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. *Advances in Neural Information Processing Systems*, 36:59662–59688.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. 2024. OlympiadBench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3828–3850, Bangkok, Thailand. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
 2021a. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.

- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Chongwen Huang, Sha Hu, George C. Alexandropoulos, Alessio Zappone, Chau Yuen, Rui Zhang, Marco Di Renzo, and Merouane Debbah. 2020. Holographic MIMO surfaces for 6G wireless networks: Opportunities, challenges, and trends. *IEEE Wireless Communications*, 27(5):118–125.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2024. Mlagentbench: Evaluating language agents on machine learning experimentation. In *Forty-first International Conference on Machine Learning*.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. arXiv preprint arXiv:2310.06825.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*.
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. 2016. Mawps: A math word problem repository. In *Proceedings of the 2016 conference of the north american chapter of the association for computational linguistics: human language technologies*, pages 1152–1157.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. 2022. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843– 3857.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 158–167, Vancouver, Canada. Association for Computational Linguistics.
- Mengbing Liu, Xin Li, Boyu Ning, Chongwen Huang, Sumei Sun, and Chau Yuen. 2022. Deep learningbased channel estimation for double-ris aided massive mimo system. *IEEE Wireless Communications Letters*, 12(1):70–74.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521.

Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, MING GONG, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie LIU. 2021. CodeXGLUE: A machine learning benchmark dataset for code understanding and generation. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.

810

811

812

814

819

821 822

823

824

825

826

827

828

831

834

837

838

839 840

841

843

852

854

855

857

858

862

- Ali Maatouk, Kenny Chirino Ampudia, Rex Ying, and Leandros Tassiulas. 2024. Tele-LLMs: A series of specialized large language models for telecommunications. *arXiv preprint arXiv:2409.05314*.
- Ali Maatouk, Fadhel Ayed, Nicola Piovesan, Antonio De Domenico, Merouane Debbah, and Zhi-Quan Luo. 2023. TeleQnA: A benchmark dataset to assess large language models telecommunications knowledge. *arXiv preprint arXiv:2310.15051*.
- Seyed Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2025. GSM-symbolic: Understanding the limitations of mathematical reasoning in large language models. In *The Thirteenth International Conference on Learning Representations*.
- OpenAI. 2024. Hello gpt-40. Accessed: 2024-10-01.
- OpenAI. 2024. Learning to reason with llms. Accessed: 2025-09-12.
- Jiawei Shao, Jingwen Tong, Qiong Wu, Wei Guo, Zijian Li, Zehong Lin, and Jun Zhang. 2024. WirelessIlm: Empowering large language models towards wireless intelligence. *arXiv preprint arXiv:2405.17053*.
- Qingjiang Shi, Meisam Razaviyayn, Zhi-Quan Luo, and Chen He. 2011. An iteratively weighted mmse approach to distributed sum-utility maximization for a MIMO interfering broadcast channel. *IEEE Transactions on Signal Processing*, 59(9):4331–4340.
- Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. 2025. Can LLMs generate novel research ideas? a largescale human study with 100+ NLP researchers. In *The Thirteenth International Conference on Learning Representations*.
- Q.H. Spencer, A.L. Swindlehurst, and M. Haardt. 2004. Zero-forcing methods for downlink spatial multiplexing in multiuser MIMO channels. *IEEE Transactions on Signal Processing*, 52(2):461–471.
- Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. 2024. Scieval: A multi-level large language model evaluation benchmark for scientific research. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 38, pages 19053–19061.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*. 863

864

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

- 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*.
- Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. 2024. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019a. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In International Conference on Learning Representations.
- Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2024. Scibench: Evaluating college-level scientific problem-solving abilities of large language models.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *Preprint*, arXiv:2409.12122.
- Haifan Yin, David Gesbert, Miltiades Filippou, and Yingzhuang Liu. 2013. A coordinated approach to channel estimation in large-scale multiple-antenna systems. *IEEE Journal on Selected Areas in Communications*, 31(2):264–273.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. Metamath: Bootstrap your own mathematical questions for large language models. In *The Twelfth International Conference on Learning Representations*.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.

- Cedegao E Zhang, Katherine M Collins, Adrian Weller, and Joshua B Tenenbaum. 2023. Ai for mathematics: A cognitive science perspective. *arXiv preprint arXiv:2310.13021*.
 - Jiakang Zheng, Jiayi Zhang, Hongyang Du, Dusit Niyato, Bo Ai, Mérouane Debbah, and Khaled B. Letaief. 2024. Mobile cell-free massive MIMO: Challenges, solutions, and future directions. *IEEE Wireless Communications*, 31(3):140–147.
 - Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. 2022. miniF2F: a cross-system benchmark for formal olympiad-level mathematics. In *International Conference on Learning Representations*.
 - Hang Zou, Qiyang Zhao, Yu Tian, Lina Bariah, Faouzi Bader, Thierry Lestable, and Merouane Debbah. 2024. TelecomGPT: A framework to build telecomspecfic large language models. *arXiv preprint arXiv:2407.09424*.

A Dataset Details

920

921

923

924

925

927

928

931

932

934

939

967

969

A.1 Topics and Papers Selection

Our study addresses modern wireless communica-940 tion challenges by selecting topics based on three 941 942 key criteria that ensure both the academic rigor and practical relevance of our work. First, we target topics that have undergone peer review and have been accepted by prestigious journals such as IEEE 945 Transactions on Wireless Communications (TWC), 946 IEEE Transactions on Communications (TCOM), 947 and IEEE Journal on Selected Areas in Communi-948 cations (JSAC), as well as by top-tier conferences 949 including IEEE International Conference on Communications(ICC) and IEEE Global Communica-951 tions Conference (Globalcom), and for which corre-952 sponding arXiv versions are available. Second, we focus on communication system modeling that en-954 tails complex, multi-step mathematical derivations. These derivations are designed to closely mimic the challenges encountered in real-world wireless 957 communication scenarios, capturing the intricate in-958 terplay between theoretical constructs and practical system constraints. Third, we ensure topic diversity 960 by covering a wide range of wireless communica-962 tion scenarios and problem domains. Specifically, our study encompasses seven major communica-963 tion scenarios and six key problem areas, including 964 interference management, spectrum optimization, 965 network coding, and energy efficiency. 966

A.2 ArXiv Data Processing

Our data processing pipeline is similar with (Maatouk et al., 2024). First begins with the removal of all comments from the LaTeX files using Google's 970 arXiv LaTeX Cleaner¹. We then parse the LaTeX 971 source to extract the core technical content by sep-972 arating the main text and mathematical expressions 973 from non-essential elements such as comments, 974 figures, and tables. For submissions comprising 975 multiple files linked via \input commands, we uti-976 lize the latexpand $tool^2$ to flatten the document 977 into a single file, ensuring all dependencies are re-978 solved. To address the variability introduced by 979 author-defined macros (e.g., via \newcommand or 980 \def), we automatically expand these using the 981 de-macro³, replacing custom macros with their 982 full definitions and normalizing all mathematical 983 expressions to a consistent LaTeX format. Non-984 informative content such as acknowledgments and 985 extensive bibliographies are removed to focus on 986 technical material and to ensure anonymity in 987 dataset construction by removing all author infor-988 mation from the articles. 989

B Prompt Templates

For clarity and reproducibility, we provide examples of our prompt templates. Figure 10 shows a template for a paper summary prompt, Figure 11 illustrates a question generation prompt, and Figures 12 and Figure 14 present templates for requesting LLMs to answer multiple-choice and fill-in-the-blank questions, respectively.

990

991

992

993

994

995

996

997

998

999

1001

1003

1004

C Model Configurations and Hyperparameters

In our WirelessMathBench experiment, a total of 16 models were tested. All model tests followed the same template and set default parameters, and all results are the results of a single run.

C.1 Closed-source Models

For models such as OpenAI-01, GPT-40, GPT-4,1005GPT-3.5-turbo, Gemini-2.0-flash, Gemini-1.5-pro,
and Gemini-1.5-flash, we utilize their official API1007interfaces. These models are invoked via their re-
spective API endpoints with standardized default
parameters to ensure consistency and reproducibil-
ity across all experiments.1005

³https://ctan.org/pkg/de-macro

¹https://github.com/google-research/ arxiv-latex-cleaner

²https://ctan.org/pkg/latexpand

1013

- 1014
- 1015
- 1016
- 1017
- 1018

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1019

1022

C.2 Open-source Models

Our local experiments employ several open-source models deployed across different environments:

- AliyunCloud Deployment⁴: Qwen2.5-Math-72B and DeepSeek-V3 are deployed on the AliyunCloud platform.
 - NVIDIA NIM Deployment⁵: DeepSeek-R1 and LLaMA-3.3-70B are run on NVIDIA NIM cloud platform.
- HuggingFace Transformers⁶: Other models-including Qwen2.5-Math-7B, LLaMA-3-8B-Tele, LLaMA-3-8B, and Mistral-7B-are run on local Nvidia A6000s, using the HuggingFace Transformers library to load pre-trained models.

Example Output D

To illustrate the diversity and complexity of the tasks within WirelessMathBench, we now present representative example outputs generated by various models. Figures 5, Figure 6, Figure 7, Figure 8, and Figure 9 showcase different task types along with the responses provided by leading models (OpenAI-o1 and DeepSeek-R1).

Disclosure 1034

In the process of writing this paper, we partially 1035 utilized ChatGPT as a language polishing tool to 1036 improve the clarity and quality of the text. However, all research ideas, data analyses, and conclu-1038 sions were independently conceived and confirmed 1039 by the authors. 1040

⁴https://bailian.console.aliyun.com/

⁵https://build.nvidia.com/

⁶https://huggingface.co/docs/transformers/en/ index

Q_id:liu2023detecting_q3_mc.



Figure 5: An example question and the corresponding output from LLMs for a multiple-choice task.

Backgroud	In an IRS-assisted MU-MIMO system, the base station receives signals through both direct and IRS-reflected paths. The observation matrix X_k combines the channel matrix H_k , IRS phase-shift matrix P , and additive noise Z_k . Here, \$M\$ is BS antennas, N IRS elements, and C pilot symbols.
Question	What noise term completes the received signal equation?
Equation	$oldsymbol{X}_k = oldsymbol{H}_k oldsymbol{P} + [MASK]$
Correct Answer	$oxed{Z}_k$
OpenAI-o1	$oxed{Z}_k$
{	The additive noise term (\boldsymbol{Z}_k) completes the received signal equation by r random disturbances in the system.", ad \boldsymbol{Z}_k "
DeepSeek-R1	
{ "explanation": " matrix \$\boldsymbol background cc \$\boldsymbol is the additive noise is added "answer": "\boxo	$\label{eq:linearly} \begin{split} & Z_k \\ \label{eq:linearly} In the given IRS-assisted MU-MIMO system, the received signal equation combines the channel symbol{H}_k$, IRS phase-shift matrix \boldsymbol{P}, and additive noise $\{Z\}_k$. The additive noise term \boldsymbol{Z}_k is explicitly mentioned in the ontext as part of the observation matrix \boldsymbol{X}_k. The equation structure $\X_k = \boldsymbol{H}_k \boldsymbol{P} + [MASK]$ directly implies that the missing term noise \boldsymbol{Z}_k, which aligns with standard wireless communication models where I linearly to the signal component.", ed{\boldsymbol{Z}_k}" \label{eq:linearly}$

Figure 6: An example of a level 1 progressively masked fill-in-the-blank task. The figure shows both the question and the LLM's output as it begins to reconstruct the missing parts of the equation.

	41_111111_ <i>2</i> .
Backgroud	In RIS-aided ISAC systems, the effective channel combines direct BS-user links and RIS-reflected paths. The RIS phase matrix Θ enables configurable signal reflections. Channel matrices maintain $N_U \times N_B$ dimensions through proper conjugation.
Question	What matrix configures RIS phase shifts?
Equation	$\mathbf{H}_{ ext{c},k} = \mathbf{H}_{ ext{BU},k} + [MASK] \mathbf{\Theta}[MASK]$
Correct Answer	$egin{array}{c} \mathbf{H}_{\mathrm{RU},k} \end{array} ext{ and } \mathbf{H}_{\mathrm{BR}}^{H} \end{bmatrix}$
OpenAI-o1	$egin{array}{c} \mathbf{H}_{\mathrm{RU},k} \end{array}$ and $egin{array}{c} \mathbf{H}_{\mathrm{BR}} \end{array}$
{ "explanation": "T (\mathbf{H}_{ "answer": "\boxed }	he two missing channel matrices are \(\mathbf{H}_{\text{RU},k}) and \ \text{BR}}).", d{\mathbf{H}_{\text{RU},k}},boxed{\mathbf{H}_{\text{BR}}}"
DeepSeek-R1	$\mathbf{H}_{\mathrm{RU},k}$ and \mathbf{H}_{BR}
{	n RIS-aided systems, the effective channel combines the direct BS-user path (H_{BU,k}) and ed path. The reflected path involves the channel from BS to RIS (H_{BR}) and RIS to user <i>i</i> th the RIS phase matrix Θ in between. Matrix multiplication order ensures dimensions align: $U \times N_R$ · Θ (N_R × N_R) · H_{BR} (N_R × N_B) results in N_U × N_B, matching us, the missing terms are H_{RU,k} and H_{BR}.", d{\mathbf{H}_{\text{BR}}}"

Figure 7: An example of a level 2 progressively masked fill-in-the-blank task.

Backgroud	Satellite line-of-sight channels use planar antenna arrays with spatial correlation. The LoS component $\overline{\mathbf{g}}_k$ depends on Rician factor κ_k , path-loss β_k , and user angles θ_k (elevation), ω_k (azimuth). Array response vector $\pm b\{ell\}$ uses wavelength λ and element positions \mathbf{c}_n .
Question	Provide scaling factor and first element's phase
Equation	$bar \mathbf{g}_k = [MASK] \Big[e^{j[MASK]}, \dots, e^{j \boldsymbol{\ell}(heta_k, \omega_k)^T \mathbf{c}_N} \Big]^T$
Correct Answer	$\boxed{\sqrt{\frac{\kappa_k \beta_k}{\kappa_k + 1}}} \text{and} \boxed{\boldsymbol{\ell}(\theta_k, \omega_k)^T \mathbf{c}_1}$
OpenAI-o1	$\boxed{\sqrt{rac{\kappa_keta_k}{\kappa_k+1}}} \hspace{0.5cm} ext{and}\hspace{0.5cm} \boxed{0}$
{ "explanation": " due to the Rid zero, so the e "answer": "\box }	The scaling factor for the LoS component is $\left\{\frac{\frac{k}{kappa_k} + 1}\right\}$ cian factor and path-loss. The first element's phase is zero because the array factor at the origin is xponential term becomes $e^{j0}=1."$, $ed\left\{\frac{\sqrt{kappa_k} + 1}\right\}$, $boxed\{0\}$ "
DeepSeek-R1	$\boxed{\sqrt{\frac{\kappa_k \beta_k}{\kappa_k + 1}}} \text{and} \boxed{\boldsymbol{\ell}(\theta_k, \omega_k)^T \mathbf{c}_1}$
{ "explanation": " \$\beta_k\$ and element's pha \omega_k)\$ v "answer": "\box \omega k)^T	The scaling factor for the LoS component $\bar{g}_k \$ combines the path-loss d Rician factor $\bar{g}_k \$ as $\bar{g}_k \$. The first se is determined by the inner product of the array response vector function $\bar{g}_k \$, with the first antenna element's position $\bar{g}_k \$,

Figure 8: An example of a level 3 progressively masked fill-in-the-blank task. Notice that both OpenAI-o1 and DeepSeek-R1 encountered difficulties in correctly interpreting the conjugate condition required by the task.

	In a stacked intelligent metasurface (SIM) with L layers and N meta-atoms per layer, each layer applies a diagonal phase-shift matrix to incoming signals. Denote the matrix of the l -th layer by $\mathbf{\Phi}^{l}$. The diagonal entries contain the
Backgroud	phase terms $e^{j\theta_n^l}$, controlling how each meta-atom shifts the wave\u2019s phase.","question_text":"Write the full expression for the phase-shift matrix of layer <i>l</i> .
Question	"Write the full expression for the phase-shift matrix of layer <i>l</i> .
Equation	$\mathbf{\Phi}^l = [MASK]$
Correct Answer	$\operatorname{diag}\!\left(e^{j heta_1^l},e^{j heta_2^l},\ldots,e^{j heta_N^l} ight)$
OpenAI-o1	$egin{aligned} & \operatorname{diag}\!\left(e^{j heta_1^l},e^{j heta_2^l},\ldots,e^{j heta_N^l} ight) \end{aligned}$
{ "explanation": ' e^{j\theta_n^l} \) function and follo \Bigr) \), commas "answer": "\box }	The phase-shift matrix $(\langle boldsymbol \{ Phi \}^1 \rangle)$ is a diagonal matrix with elements $(\langle along the diagonal for each meta-atom (n = 1, 2, \langle dotsc, N \rangle). Using the \langle (mathrm \{ diag \} \rangle) wing the LaTeX protocols: functions in roman font, scaled parentheses via \langle (Bigl(\rangle) and \langle (with small spaces via \langle (, \rangle), and braced exponents \langle (e^{\{ (\dots,)\}} \rangle), the expression becomes:", ed{\langle mathrm \{ diag \} \langle Bigl(e^{j} theta^{l_1}, e^{j} theta^{l_2}, \langle dotsc, \rangle, e^{j} theta^{l_1}N \rangle \langle Bigr) \rangle"$
DeepSeek-R1	$\operatorname{diag}\!\left(e^{j heta_1^l},e^{j heta_2^l},\ldots,e^{j heta_N^l} ight)$
{ "explanation": ' corresponds to the scaled parentheses	The phase-shift matrix for layer l in a SIM is a diagonal matrix where each diagonal entry phase shift $e^{j\theta_n^l}$ applied by the <i>n</i> -th meta-atom in the <i>l</i> -th layer. Using the diag function with and properly formatted arguments, the matrix is constructed as diag $\left(e^{j\theta_1^l}, e^{j\theta_2^l}, \ldots, e^{j\theta_N^l}\right)$.",

Figure 9: An example of a Full Equation Completion (FEC) task.

Paper Summary Prompt Template.

Task: Act as an expert in wireless communications and mathematical modeling. You will be provided with a research paper in the field of wireless communications. Your goal is to extract and summarize the mathematical system modeling contained in the paper. <<STRUCTURE REQUIREMENTS>> 1. **Model Extraction**: a) Identify ALL system equations with context b) For each equation: i) List ALL variables with units/dimensions

- ii) Specify underlying assumptions
- iii) Note domain restrictions

2. **Summary Organization**:

\paragraph {Background} (2-3 sentences contextualizing the model) \paragraph {Key Assumptions} (bullet points with \$\bullet\$) \paragraph {Parameter Definitions} (table-like structure using \quad) \paragraph {Core Equations} (numbered with original labels if available)

- 3. **Equation Formatting**:
- Vectors: \boldsymbol {v}
- Matrices: \mathbf{M}
- Operators: \mathrm {diag}, \mathrm {tr}
- Complex numbers: j for imaginary unit
- Dirac notation: \delta(\cdot)
- Alignment: \begin {align*}...\end {align*}
- <<CONTENT GUIDELINES>>
- **Variable Explanations**:
- For each symbol: \$\theta\$ \quad (Type: Phase shift; Domain: [0,2π); Unit: rad)
- Matrix dimensions: $\operatorname{H} \in \mathbb{C}^{N \times M}$
- Distinguish similar symbols: $h_{ij} \ vs h_{i}^{(j)}$
- ***Model Validation**:
- Verify dimensional consistency
- Check boundary conditions
- Confirm parameter unit homogeneity
- **Special Notation**:
- Estimated quantities: \hat{x}
- Time derivatives: $\det{\det{v}}$
- Optimal values: \$\mathbf{\Theta}^\ast\$

<<EXAMPLE OUTPUT>>>

\paragraph {Background}
Consider an RIS-assisted mmWave system with \$K\$ single-antenna users
communicating via \$N\$ RIS elements...
<Key Assumptions>
\textbf{1} Effective Channel Representation.}\\
The received signal at the BS from user \(k\) combines the single-reflection and
double-reflection links.....
The equivalent channel from user \(k\) to the BS is
\begin {equation} \label{eq:EffectiveChannel}
........
\end{equation}
\end{equation}
<Parameter Definitions> such as \$\boldsymbol{h}_{k1} \in \mathbb{C}^{N \times
1}\$ is the channel from user \(k\) to \(\mathcal{RIS}_1\).

\paragraph{System model with equations.}



Question Generation Prompt Template.

Task: Act as an expert in wireless communications and mathematical modeling You will generate exam-style questions from research paper summaries with the following strict requirements:

<<STRUCTURE REOUIREMENTS>>

1. **Per Equation Processing*

a) Identify ALL system model equations, explain each symbol in the equation in background

- b) For EACH equation:i) Mask the RHS with [MASK]

 Mask the RHS with [MASK]
 Generate 1 MCQ with 4 plausible options
 Create 4 progressive fill-in-the-blank subquestions with:
 25%, 50%, 75%, and 100% key symbols masked
 Variable masking in different positions
 iv) Ensure each subquestions can be answered independently and full self-timeted in external comparation containment (no external references)

- 2. **Question Components**
- **Question Components**:
 For MCQs:
 * textbf{Background}: Contextual info in 3-5 lines
 * textbf{Background}: Contextual info in 3-5 lines
 * textbf{Equation}: Masked equation in display math mode
 * textbf{Question}: Explicitly ask to replace [MASK]
 * textbf{Options}: 4 LaTeX-formatted choices (A)-(D)
 * textbf{Options}: 4 LaTeX-formatted choice
- * \textbf{Answer}: Detailed derivation walkthrough For Fill-in-the-blanks:

- * \textbf{Masked Equation}: Progressive symbol replacement * \textbf{Question}: Specific term request
- * \textbf{Answer}: Complete equation with \boxed{} for solutions
- <<CONTENT GUIDELINES>> **Distractor Design**:

- Below are some common error types:
 Matrix dimension mismatches
- 2) Incorrect sequence $(\theta_1 \text{ vs } \theta_2)$
- 3) Missing diag() operators
- 4) Channel matrix transposition errors
- 5) Incorrect matrix multiplication
- 6) Incorrect vector/matrix dimensions 7) Random mistakes
- **Mathematical Rigor**
- Use \mathrm{diag}() for diagonalization
- Replace symbolic writing in the problem appropriately, but ensure correctness
 Verify matrix multiplication compatibility
- **Pedagogical Elements**:
 - Add \underbrace annotations in answers Use \$;\$; for proper equation spacing
 - Include \text{} descriptions for terms
- **Masking Strategy**:
 Mask variables not operators
- Progressively increase masked terms:
 1) 25%: Single critical variable
- 2) 50%: Two interdependent terms
 3) 75%: Multiple components
- 4) 100%: Full equation recall
- **Blank Positioning**:
- 50% mask: Key parameter (e.g., θ values)
 50% mask: Combined terms (e.g., N_jdiag(h_kj))
 75% mask: Structural components
 100% mack: Entrie DUC
- 100% mask: Entire RHS
- <<FORMATTING RULES>>
- Strict LaTeX compliance:
 Equations in \$\$ \$\$/equation* environments
- \boldsymbol for vectors/matrices
 \mathrm for operators (diag, etc)
- \quad spacing between equation terms
- · Section headers with: \textbf{\Large Question X (System Component)} \vspace{-0.5em} % Compact vertical spacing

<<EXAMPLE TEMPLATE>>>

Figure 11: This prompt template converts paper summaries into detailed question-answer pairs.

Multiple-Choice Question Prompt Template.

following multiple-o	choice question.
Background Cont {background}	ext
Multiple-Choice ({question_text}	Question
Relevant Equation {equation or "No eq	n uation provided"}
{options_str}	
Formatting Requi {latex_core_rules} - Explanations requi - Maintain consisten **Response Instruct 1. Select ONE letter 2. Provide technical 3. Present your final explanation	rements ring math must use for inline equations it notation with question context tions** choice (A-D) reasoning in the explanation field answer in a JSON format. This should include: answer and
Your output should blocks tagged with ' `json {{ "explanation": "< "answer": "C", }}	be formatted as a JSON object enclosed in Markdown code json'. For example: explanation>"

Figure 12: This template is designed for answering multiple-choice questions. The model is guided to choose the correct mathematical expression from a set of closely related options.

Fill-in-the-Blank and Full Equation Completion Question Prompt Template. You are a domain expert in Wireless Communication. Please answer the following fill-in-the-blank question. **Background Context** {background} **Problem Statement (Blanks marked with [MASK])** {question text} **Equation to Complete** {equation or "No equation provided"} **Strict LaTeX Protocol** {latex core rules} 5. Each [MASK] requires a separate \boxed { {... } } 6. Final answer line format: The final answer is $boxed{\{...\}}, boxed{\{...\}},...$ **Submission Requirements** 1. Answer box number corresponds to the [MASK] number 2. Use EXACT formatting from the equation/question 3. No natural language in boxed answers 4. Technical explanation in the 'explanation' field Your output should be formatted as a JSON object enclosed in Markdown code blocks tagged with 'json'. For example: `json {{ "explanation": "<explanation>", "answer": "\boxed{<answer2>}, \boxed{<answer3>}", }}

Figure 13: This prompt template is used for fill-in-the-blank and full equation completion tasks. It directs the model to reconstruct missing parts of equations by using contextual cues and domain knowledge, simulating the process of step-by-step derivation.

LLM Evaluation Answer Prompt Template.

You are an expert in wireless communications and mathematical modeling. Your task is to evaluate a student's answer against the correct answer. Follow these evaluation criteria strictly:

1. ******Mathematical Equivalence:****** Check if the student's answer is mathematically equivalent to the correct answer.

2. ******Answer Format:****** The student's answer should be in the same format as the correct answer.

3. ******All blanks should be filled.****** All the blanks should be filled with the correct answer.

4. **Scoring:** Output a score of "1" if the student's answer is correct (i.e., mathematically equivalent) or "0" if it is not.

Please provide your final output as a JSON object (enclosed in Markdown code blocks tagged with 'json') with the following format:

```json {{ "score": "1" ;;; or ```json {{ "score": "0" <u>}}</u> Here is the background: {background} Here is the question: {question} Here is the question: {question} Here is the correct answer: {true\_answer} Here is the student's answer: {student\_answer}

Evaluate the student's answer based solely on the above information and output only the JSON object with the score.

Figure 14: This evaluation prompt template standardizes the process of assessing model-generated answers. It ensures that responses are judged consistently based on their correctness, completeness, and adherence to the required domain-specific reasoning.