SalaMAnder: Shapley-based Mathematical Expression Attribution and Metric for Chain-of-Thought Reasoning

Anonymous ACL submission

Abstract

While chain-of-thought (CoT) prompting has 001 002 demonstrated remarkable efficacy in enhancing the reasoning capacities of large language models (LLMs) for mathematical problemsolving, the mechanistic foundations underlying these improvements remain inadequately characterized and lack systematic theoretical investigation. In this paper, we present Sala-MAnder (Shapley-based Mathematical Expression Attribution and Metric), a theoreti-011 cally grounded methodology as well as a math-012 ematically rigorous evaluation metric for quantifying component-level contributions in CoT reasoning. Specifically, we leverage Shapley value for mathematical expression attribution and develop an efficient stratified sampling al-017 gorithm that significantly reduces the computational complexity. Besides, we develop the CoSP (Cardinality of Shapley Positives) metric through covariance analysis. Comprehensive validation across multiple LLM models and diverse mathematical benchmarks demonstrate that the CoSP metric within our SalaMAnder framework exhibits a robust monotonic correlation with model performance. This correlation not only provides theoretical explanations for 027 the empirical success of existing CoT but also establishes mathematically rigorous principles for prompt construction optimization. Finally, the analytical capabilities of SalaMAnder is further substantiated by performance improvements achieved through targeted refinement of low-CoSP components, demonstrating both the explanatory power and practical utility in understanding and enhancing CoT reasoning.

1 Introduction

036

042

The emergency of chain-of-thought (CoT) reasoning has propelled large language models (LLMs) to achieve human-level performance in complex mathematical reasoning tasks, from arithmetic problem solving to theorem proving. Despite the empirical advances, the field confronts a fundamental scientific challenge: current understanding of why specific reasoning steps lead to correct solutions remains trapped in a cycle of heuristic speculation(Wang et al., 2023; Chen et al., 2024; Wang et al., 2022; Li et al., 2024; Jin et al., 2024; Pfau et al., 2024) or labor-intensive verification(Serrano and Smith, 2019; Bastings and Filippova, 2020; Madsen et al., 2022; Siddiqui et al., 2024), lacking systematic theoretical investigation. We reveal that this issue stems from two fundamental limitations in existing interpretation methodologies. For one thing, current approaches predominantly depend on heuristic-driven engineering practices, where practitioners optimize CoT demonstrations through ad hoc trial-and-error adjustments or case-specific manual inspections. This reliance on empirical intuition rather than systematic analysis yields explanations that lack both mathematical rigor and generalizable insights. For another thing, while approaches such as exact Shapley value computation (Shapley, 1953; Weber, 1988) provide mathematical rigor, their exponential complexity renders them impractical for real-world applications. For example, a single Chain-of-Thought demonstration containing just 30 components necessitates over 1 billion evaluations. The empirical fragility in heuristic methods and computational infeasibility in rigorous approaches have significantly impeded the development of scalable, principled frameworks for systematic CoT analysis and optimization.

043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

In this paper, we propose a unified framework, namely **SalaMAnder** (Shapley-based Mathematical Expression Attribution and Metric), that introduces two key innovations for efficient and semantically coherent CoT analysis. First, we establish mathematical expressions as atomic units for Shapley-based attribution, addressing the semantic fragmentation inherent in traditional tokenlevel analyses through component-level decomposition. Then, we develop a novel stratified sampling algorithm, namely SalaMA (Shapley-based

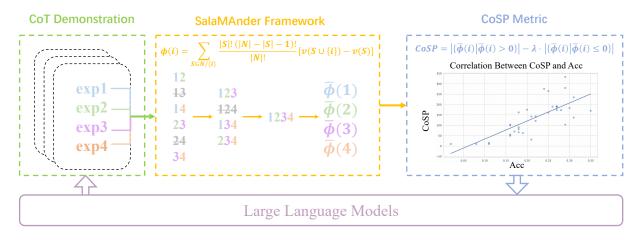


Figure 1: Workflow of the SalaMAnder Framework and CoSP Metric in CoT for LLMs. Initially, the framework proposes an efficient Shapley value algorithm to attribute the contributions of various mathematical expressions. These computed Shapley values are then utilized to derive the CoSP metric. Both theoretical derivations and extensive experiments across multiple models and datasets validate that CoSP exhibits a robust positive correlation with model inference accuracy. This correlation provides a comprehensive explanation of the underlying mechanisms driving CoT behavior in LLMs.

Mathematical Expression Attribution) that achieves exponential complexity reduction by decomposing Shapley calculations according to component order, reducing time complexity from $O(2^{n+1})$ to $O(2mn^2)$, while maintaining rigorous theoretical guarantees. Besides, we develop the **CoSP** (Cardinality of Shapley Positives) metric based on the efficient and semantical shapely estimation. The proposed CoSP metric within our SalaMAnder framework formally establishes the monotonic relationship with model performance through rigorous covariance analysis, providing mathematical guarantees for the predictive validity.

084

086

090

100

103

104

105

108

109

110

The contributions of this paper are summarized:

- We propose a unified framework, namely Sala-MAnder that establishes mathematical expressions as atomic units for Shapley-based attribution and develop a novel stratified sampling algorithm, namely SalaMA that achieves exponential complexity reduction while maintaining rigorous theoretical guarantees.
- We propose the CoSP metric within our SalaMAnder framework, which formally establishes the monotonic relationship with model performance through rigorous covariance analysis, providing mathematical guarantees for the predictive validity.

Experimentally, we first utilize SalaMAnder in fewshot learning scenarios to assess the validity of our
explanation method and metric. Then we further
evaluate the reliability of explanation results. Last

we present novel insights that not only reinforce the effectiveness of our methods but also integrate and unify previous research. 115

116

117

118

119

120

121

122

123

124

125

127

128

129

130

131

132

133

134

135

136

137

139

140

141

142

143

144

2 Related Work

CoT Methodologies CoT prompting, introduced by Wei et al. (2022), explicitly guides LLMs to generate intermediate reasoning steps, significantly improving performance on mathematical and symbolic tasks. Subsequent work expanded this paradigm through path optimization (e.g., Least-to-Most prompting decomposes problems into subquestions (Zhou et al., 2022); Progressive-Hint iteratively refines solutions (Zheng et al., 2023)), automation (e.g., Automatic CoT generates demonstrations via LLMs (Zhang et al., 2022); Symbolic CoT Distillation transfers CoT ability to smaller models (Li et al., 2023)), and hybrid approaches (e.g., CoF-CoT combines coarse-to-fine prompting for multi-domain tasks (Nguyen et al., 2023); Deductive Verification adds formal consistency checks (Ling et al., 2023)). Despite these advances, most methods rely on heuristic designs without theoretical guarantees, and their efficacy varies significantly across domains-mathematical tasks benefit more from structured CoT than openended reasoning.

Mechanistic Studies of CoT Reasoning The existing literature on CoT mechanisms unfolds through complementary empirical and theoretical lenses. Empirical studies (Wang et al., 2022; Li

et al., 2024; Jin et al., 2024; Wang et al., 2023; Pfau 145 et al., 2024; Chen et al., 2024) have explored vari-146 ous strategies to enhance the robustness, safety, and 147 structural integrity of CoT reasoning. For instance, 148 self-consistency mechanisms (Wang et al., 2022) improve the reliability of reasoning outputs by ag-150 gregating multiple reasoning paths, while efforts 151 to mitigate toxicity (Li et al., 2024) ensure safer 152 commonsense reasoning. Additionally, research on 153 step length (Jin et al., 2024), step relevance and 154 logical order (Wang et al., 2023), hidden state dy-155 namics (Pfau et al., 2024), and premise sequence 156 order (Chen et al., 2024) underscores the importance of prompt design and structural factors in 158 optimizing CoT performance. 159

3 Method

160

161

162

163

164

167

168

169

170

171

172

173

174

175

176

177

178

179

180

182

184

185

In this section, we introduce the **SalaMAnder** framework, designed to explain the mathematical reasoning mechanisms of CoT in LLMs using Shapley values. We introduce our method in three sections: an introduction to Shapley values, the **SalaMAnder** sparse computation of these values, and the **CoSP** metric for evaluating CoT reasoning contributions.

3.1 Preliminary: Shapley Values (Fair Attribution of CoT Constituents)

Shapley values, originating from cooperative game theory, offer a principled method for fairly distributing the total gains of a coalition among its individual players based on their contributions (Shapley, 1953).

Formally, consider a set of players $N = \{1, 2, ..., n\}$ and a reward function $v : 2^N \to \mathbb{R}$ that assigns a real-valued payoff to every possible coalition of players. The Shapley value $\phi_i(v)$ for player *i* is defined as:

$$\phi_{v}(i) = \sum_{S \subseteq N \setminus \{i\}} \frac{s!(n-s-1)!}{n!} \left[v(S \cup \{i\}) - v(S) \right]$$

where S is any subset of N that does not include player i, and s = |S|, n = |N| respectively denotes the number of players in subset S and set N.

We can further derive from the above expression:

186
$$\phi(i) = \frac{1}{n} \sum_{S \subseteq N \setminus \{i\}} \frac{1}{\binom{n-1}{s}} \left[v(S \cup \{i\}) - v(S) \right]$$
187
$$= \frac{1}{n} \sum_{r=0}^{n-1} \mathbb{E}_{s=r} \left[v(S \cup \{i\}) - v(S) \right]$$
188
$$= \frac{1}{n} \phi_{r+1}(i)$$
(1)

where
$$\phi_k(i) = \mathbb{E}_{s=r} [v(S \cup \{i\}) - v(S)]$$
 denotes
the $(r+1)$ th order shapley value of component *i*.

189

190

191

192

193

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

228

229

230

231

232

233

234

235

236

237

238

Researchers have proven that the Shapley value is a unique unbiased method to fairly allocate overall reward to each player with four properties: linearity, dummy, symmetry, and efficiency (Weber, 1988). For simplicity, we use $\phi(i)$ by ignoring the superscript of $\phi_v(i)$ in the following manuscript without causing ambiguity.

In our framework, each component of the CoT, such as individual mathematical expressions or a single word, is treated as a player in the cooperative game. The reward function v(S) corresponds to a performance metric of the LLM (e.g., correctness, or inference logits) when only the components in subset S are included in the CoT. Consequently, the Shapley value $\phi(i)$ quantifies the average marginal contribution of each component to the overall reasoning performance across all possible subsets of components.

3.2 SalaMA: Efficient Sparse Shapley Computation for CoT Components

Although calculating exact Shapley values for each component presents significant computational challenges, the exponential growth in the number of possible subsets with respect to the number of components renders exact computation infeasible for practical applications. To address the limitation, we propose SalaMA (Shapley-based Mathematical Expression Attribution) mechanism, an efficient algorithm designed to approximate Shapley values with high accuracy while substantially reducing computational overhead.

The Players We define each player in the game, i.e. each component in the demonstration as a mathematical expression rather than individual words or tokens. This decision is motivated by the observation that single words or tokens can vary in meaning across different contexts, making their attribution inconsistent and less meaningful. Mathematical expressions, in contrast, maintain their semantic integrity across diverse reasoning scenarios, providing a more stable and universally applicable unit for analysis. Additionally, aggregating tokens into coherent mathematical expressions significantly reduces the number of components, thereby mitigating the computational complexity associated with Shapley value calculations. This aggregation not only enhances computational efficiency but also ensures that the attribution analysis

241 242

243

244

247

262

264

265

267

263

268

269

271

273

274

275

278

279

282

remains interpretable and relevant to the model's 239 problem-solving mechanisms.

The Reward Function We adopt a reward function that combines the model's prediction confidence logits with the correctness of the prediction, formulated as

$$\begin{cases} v(S) = \left(\frac{1}{L} \sum_{\ell=1}^{L} \log p_{\theta}(y_{\ell}|S)\right) \cdot \mathbb{I}(y_{\text{pred}}(S) = y^{*} \\ y_{\text{pred}}(S) = \bigoplus_{\ell=1}^{L} y_{\ell}(S) \end{cases}$$

$$\tag{2}$$

where $\frac{1}{L}\sum_{\ell=1}^L \log p_{\theta}(y_{\ell}|S)$ represents the average confidence score of the model's prediction by averaging the logits associated with the result tokens generated when including component subset $S, \mathbb{I}(\cdot)$ is a binary indicator, and \bigoplus indicates the string concatenation operation.

This formulation ensures that the value function directly reflects the impact of each component on the model's performance, addressing the limitations of alternative metrics such as attention or saliency scores or binary correctness. Attention or saliency scores do not provide a direct attribution to the final outcome and can be complex to interpret (Serrano and Smith, 2019; Bastings and Filippova, 2020; Madsen et al., 2022; Siddiqui et al., 2024), while a binary correctness metric lacks the sensitivity needed to capture nuanced contributions. By integrating confidence logits with correctness, reward function balances sensitivity and direct attribution, facilitating a more accurate and interpretable estimation of each component's contribution.

Efficient Shapley Computation Algorithm The proposed algorithm systematically approximates the Shapley values for CoT components through a structured algorithmic workflow. In exact Shapley value computation, for each component *i*, it is necessary to evaluate $v(S \cup \{i\}) - v(S)$ across all subsets $S \subseteq N\{i\}$, leading to a computational complexity of $O(2^{n+1})$, where n is the number of components. This exponential complexity becomes prohibitively expensive as the number of components increases. To mitigate this, SalaMA reduces the number of necessary inferences by employing a stratified sampling approach based on the order of Shapley values.

> Specifically, the SalaMA mechanism decomposes the Shapley value calculation by order. For

an r-th order Shapley value ϕ_r , SalaMA randomly samples r-1 other mathematical expressions from the set $N/\{i\}$. The number of such samples is denoted by sp, with a maximum limit of m, indicating $sp = \min(m, \binom{n-1}{r-1})$. In the original demonstration, aside from the mathematical expressions, other components (referred to as the "whiteboard") are always present and remain constant across dif-) ferent subsets.

283

284

285

289

291

292

293

294

295

296

298

299

301

302

306

307

308

309

310

311

312

313

314

315

316

317

318

319

321

322

323

324

325

326

During inference, for each sampled subset S of size r-1, SalaMA constructs two distinct demonstrations: one containing $S \cup \{i\}$ combined with the whiteboard, and another containing S alone with the whiteboard. These demonstrations are then fed into the model to obtain the corresponding reward functions $v(S \cup \{i\})$ and v(S), respectively. By iterating over multiple orders and different samples within each order, SalaMA aggregates the marginal contributions across various subset configurations. The approximated shapley value can be derived from Eq. (1):

$$\phi(i) = \frac{1}{n} \sum_{r=0}^{n-1} \mathbb{E}_{s=r}[v(S \cup \{i\} - v(S))]$$
304

$$= \frac{1}{n} \sum_{r=0}^{n-1} \frac{1}{m} \sum_{t=1}^{m} [v(S_t^r \cup \{i\} - v(S_t^r))] \quad (3)$$

To further enhance computational efficiency, SalaMA maintains a hash table \mathcal{H} to store and retrieve the results of previously computed subsets S. This caching mechanism prevents redundant inferences by ensuring that once a subset S has been evaluated, its corresponding v(S) does not need to be recomputed in future iterations. Consequently, the computational complexity of SalaMA is reduced to $O(2 \cdot sp \cdot n^2) \leq O(2mn^2)$, which is significantly lower than the exact Shapley value computation's $O(2^{n+1})$. The whole workflow is shown in Algorithm. 1.

3.3 **CoSP: Performance-Aligned Causal Explanation Rationale**

We introduce CoSP (Cardinality of Shapley Positives), a metric defined as the number of expressions within a demonstration that exhibit positive average Shapley values minus a weighted nonpositive average Shapley values across multiple experiments.

Formally, for a demonstration comprising a set

Algorithm 1: SalaMA: Sparse Shapley Value Computation

```
Function SalaMA(N, v, n, m):
      Initialize \phi[i] \leftarrow 0 \ (\forall i \in N), \ \mathcal{H} \leftarrow \emptyset;
      foreach i \in N do
            for r = 1 to n do
                  sp \leftarrow \min(m, \binom{n-1}{r-1});
                  for s = 1 to sp do
                         S \leftarrow \text{Sample}(r-1, N \setminus i);
                         v_S \leftarrow \mathsf{MemEval}(S, \mathcal{H});
                         v_{S\cup i} \leftarrow \mathsf{MemEval}(S \cup i, \mathcal{H});
                         \phi[i] \mathrel{+}= (v_{S\cup i} - v_S)/(sp \cdot n);
                  end
            end
      end
      return \phi;
Procedure MemEval(S, \mathcal{H}):
      if S \notin \mathcal{H} then
           \mathcal{H}[S] \leftarrow v(S);
```

end
return
$$\mathcal{H}[S]$$
;

327

331

332

333

335

337

338

340

341

342

344

347

348

of n expressions N, CoSP is defined as:

$$CoSP = |\{\bar{\phi}(i)|\bar{\phi}(i) > 0\}| - \lambda \cdot |\{\bar{\phi}(i)|\bar{\phi}(i) \le 0\}|$$
$$= \sum_{i=1}^{n} \mathbb{I}(\bar{\phi}(i) > 0) - \lambda \cdot \mathbb{I}(\bar{\phi}(i) \le 0)$$
$$= (1+\lambda) \sum_{i=1}^{n} \mathbb{I}(\bar{\phi}(i) > 0) - \lambda n$$

where $\phi(i)$ is the average Shapley value of the *i*-th expression, computed over *m* different problem instances tested using the same demonstration, formulated as $\bar{\phi}(i) = \frac{1}{m} \sum_{k=1}^{m} \phi^{(k)}(i)$, $\mathbb{I}(\cdot)$ is the indicator function, returning 1 if the condition inside is true and 0 otherwise, and $\lambda > 0$ is the penalty severity for the number of expressions with negative Shapley values. And we assume that during the *m* CoT reasoning precesses, for each expression *i*, there is $\phi^{(k)}(i) \sim \mathcal{N}(\mu_i, \sigma_i^2)$.

A positive average Shapley value ($\bar{\phi}(i) > 0$) indicates that the corresponding mathematical expression contributes positively to the model's reasoning performance; conversely, a non-positive one leads to negative contribution or no contribution. Therefore, CoSP comprehensively quantifies the number of expressions that actively enhance or degrade the model's efficacy in solving problems. A higher CoSP suggests that a greater subset of expressions349within the CoT is beneficial while a smaller subset350harmful, correlating with improved model performance. Specifically, we define CoSP-0 and CoSP-3511, with λ equals to 0 and 1, respectively.353

354

355

357

358

362

364

365

372

373

374

375

376

377

378

379

381

382

383

384

386

387

To substantiate the relationship between CoSP and performance, we formalize the following two theorems under specific statistical assumptions.

Theorem 1 Both CoSP-0 and CoSP-1 have positive correlation with the model performance:

$$Cov(CoSP, Perf) = (1+\lambda)(\delta_{+} - \delta_{-})\sum_{i=1}^{n} Var(X_{i})$$
35

$$Cov(Perf, CoSP-0) = (\delta_{+} - \delta_{-}) \sum_{i=1}^{n} Var(X_{i})$$
360

$$Cov(Perf, CoSP-I) = 2(\delta_{+} - \delta_{-}) \sum_{i=1}^{n} Var(X_i)$$
(4)

where the meaning of δ_+, δ_-, X_i will be explained in the proof.

Theorem 2 CoSP-0 has a positive correlation with the number of expressions n, while CoSP-1 has a negative correlation with n:

$$\mathbb{E}[CoSP_{n+1}] = (1+\lambda)\sum_{i=1}^{n+1} p_i - (n+1)\lambda$$

$$= \mathbb{E}[CoSP_n] + p_{n+1} - \lambda \qquad (5)$$

$$\mathbb{E}[CoSP \cdot \theta_{n+1}] - \mathbb{E}[CoSP \cdot \theta_n] = p_{n+1} > 0$$
$$\mathbb{E}[CoSP \cdot I_{n+1}] - \mathbb{E}[CoSP \cdot I_n] = p_{n+1} - 1 < 0$$
(6)

The proof of Theo. 1 and Theo. 2 is applied in Appendix. A.

The number of expressions n in the CoT is often indicative of the complexity or difficulty of the reasoning task. Generally, increased reasoning difficulty generally leads to better model performance (OpenAI, 2024), provided that the additional complexity is constructively leveraged. Our Theo. 2 aligns with this observation by showing that a higher number of expressions n results in a higher CoSP-0, which in turn, per Theo. 1, correlates with enhanced model performance. This consistency underscores the validity of CoSP as a metric that not only accounts for the quantity of reasoning steps but also their qualitative impact on model efficacy.

Datasets	Correlation between Metrics and Model Inference Accuracy												
	LLaMA 2				LLaMA 3				Qwen 2.5				
	CoSP-0	CoSP-1	SSV	NoE	CoSP-0	CoSP-1	SSV	NoE	CoSP-0	CoSP-1	SSV	NoE	
					j	l-shot							
GSM8K	0.76	0.65	0.32	0.76	0.70	0.18	-0.14	0.71	0.64	0.62	0.54	0.43	
MathQA	0.44	0.62	0.63	-0.08	0.37	0.28	0.19	0.10	-0.16	0.28	0.11	-0.22	
AQUA	0.40	0.46	0.44	-0.31	-0.21	0.48	0.39	-0.40	-0.63	-0.03	-0.03	-0.67	
MultiArith	0.60	0.52	0.02	0.53	0.74	0.44	0.44	0.09	0.78	0.71	0.80	-0.04	
SVAMP	0.49	0.28	0.21	0.14	0.17	0.21	0.08	-0.35	0.56	0.50	0.56	-0.32	
2-shot													
GSM8K	0.75	0.35	0.14	0.75	0.49	0.26	0.24	0.45	0.80	0.48	0.51	0.13	
MathQA	0.36	0.46	0.35	-0.11	-0.20	0.01	0.07	-0.05	-0.20	-0.14	-0.03	-0.06	
AQUA	0.56	0.51	0.48	-0.47	0.09	-0.04	-0.22	-0.50	0.22	0.52	0.55	-0.19	
MultiArith	-0.04	-0.07	-0.20	-0.31	0.82	0.39	0.58	-0.24	0.44	0.18	0.16	0.06	
SVAMP	0.23	0.05	-0.13	-0.02	0.47	0.44	-0.19	-0.17	0.69	0.61	0.53	-0.02	
4-shot													
GSM8K	0.77	0.61	0.12	0.52	0.26	0.37	-0.15	-0.20	0.80	0.58	0.52	0.31	
MathQA	0.29	-0.26	-0.46	-0.01	0.40	0.28	-0.02	-0.67	0.18	-0.33	-0.52	0.14	
AQUA	0.80	0.77	-0.10	-0.11	-0.08	0.20	0.02	-0.19	-0.31	-0.11	-0.05	-0.43	
MultiArith	0.54	0.33	0.42	0.22	0.80	0.23	-0.001	-0.47	0.67	0.51	0.24	-0.44	
SVAMP	0.63	0.31	0.22	0.61	0.10	0.07	0.36	-0.17	0.22	-0.03	-0.14	-0.13	
Average	0.51	0.37	0.16	0.14	0.33	0.25	0.11	-0.14	0.31	0.29	0.25	-0.10	

Table 1: The correlation coefficients between different metrics and model inference accuracy across multiple datasets and models of few-shot tasks. For each dataset and each model, the largest correlation is **bolded**, indicating the best interpretation method.

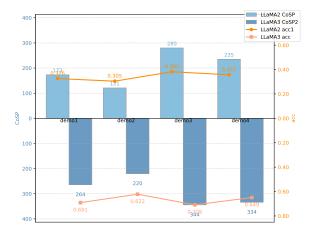


Figure 2: CoSP of LLaMA 2 and LLaMA 3.

4 Experiments

This section presents a comprehensive evaluation of the proposed SalaMAnder framework, demonstrating its applicability across various settings. Sec 4.1 describes the experimental settings, and Sec 4.2 utilizes SalaMAnder in few-shot learning scenarios to assess the validity of our explanation method and metric. In Sec 4.3, we further evaluate the reliability of explanation results. Sec 4.4 present novel insights that not only reinforce the effectiveness of our methods but also integrate and unify previous research.

4.1 Experimental Settings

To evaluate the effectiveness of the proposed SalaMA method and the CoSP metric, we conducted experiments using three foundational large language models and five representative mathematical datasets. The selected models, LLaMA-2-13Bchat (Touvron et al., 2023), LLaMA-3-8B-Instruct (Grattafiori et al., 2024), and Qwen2.5-7B-Instruct (Team, 2024) were chosen for their fundamental architectures and generalizability, as they are not overly specialized or pre-trained on extensive mathematical datasets. This ensures that our analysis of CoSP and SalaMA is broadly applicable across different model paradigms. 400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

For the datasets, we utilized GSM8K (Ouyang et al., 2022), MathQA (Amini et al., 2019), AQUA (Ling et al., 2017), MultiArith (Wang et al., 2018), and SVAMP (Patel et al., 2021). These datasets were selected for their representativeness in the mathematical question-answering domain, encompassing a range of difficulties where MathQA and AQUA are approximately equivalent and more challenging than GSM8K, which is in turn more difficult than MultiArith and SVAMP. Specifically,

Datasets	Average Correlation												
		1-sho	t		2-shot				4-shot				
	CoSP-0	CoSP-1	SSV	NoE	CoSP-0	CoSP-1	SSV	NoE	CoSP-0	CoSP-1	SSV	NoE	
GSM8K	0.70	0.48	0.24	0.63	0.68	0.36	0.33	0.44	0.61	0.52	0.16	0.21	
MathQA	0.22	0.39	0.31	-0.07	-0.01	0.11	0.13	-0.07	0.29	-0.10	-0.33	-0.18	
AQUA	-0.15	0.30	0.27	-0.46	0.29	0.33	0.27	-0.39	0.14	0.29	-0.04	-0.24	
MultiArith	0.71	0.56	0.02	0.42	0.41	0.17	0.18	-0.16	0.64	0.36	0.22	-0.23	
SVAMP	0.41	0.33	0.28	-0.18	0.46	0.37	0.07	-0.07	0.32	0.12	0.15	0.10	

Table 2: The correlation coefficients averaged among various models in few-shot tasks. For each dataset, the largest correlation is **bolded**, indicating the best interpretation method.

GSM8K consists of grade-school level math problems, MathQA includes complex multi-step reasoning questions, AQUA focuses on arithmetic and algebraic tasks, MultiArith provides multi-step arithmetic word problems, and SVAMP introduces adversarial variations to traditional arithmetic problems. This selection ensures comprehensive coverage of various aspects and complexities inherent in mathematical QA tasks.

424

425

426

427

428

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

4.2 Attribution Validity: CoSP Metric Verification in Few-Shot Learning

To evaluate the practical applicability of the proposed SalaMA method and the CoSP metric, we applied them to few-shot learning scenarios across multiple mathematical datasets and foundational language models to assess the correlation between CoSP and model performance (accuracy), thereby validating the effectiveness of our framework.

We meticulously constructed demonstrations to ensure a uniform distribution of mathematical expressions. Specifically, for one-shot learning tasks, we constructed demonstrations by selecting 35 question-answer (Q-A) pairs from the training sets of the GSM8K, MathQA, and AQUA datasets. Because the MultiArith and SVAMP datasets include answers composed solely of single mathematical expressions, we instead selected 35 Q-A pairs from the GSM8K dataset to serve as demonstrations. These one-shot demonstrations were evenly distributed, with five Q-A pairs each containing between one and seven mathematical expressions. For 2-shot demonstrations, the total number of expressions ranged from 2 to 10, resulting in 14 unique demonstrations by accounting for multiple combinations where applicable (e.g., a total of 6 expressions could be achieved by combinations 2+4 or 3+3). 4-shot demonstrations contained 4-16 total expressions, with one unique combination retained per expression count to minimize computation, producing 13 distinct demonstration sets. This methodology ensured that both one-shot and fewshot demonstrations maintained a balanced and uniform distribution of mathematical expressions, thereby isolating the effect of expression quantity on model performance. 463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

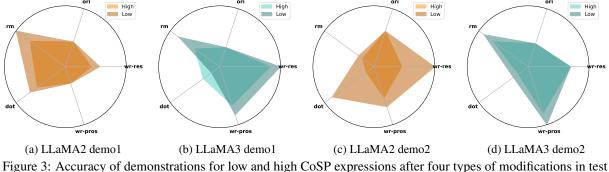
We then utilize the proposed SalaMA method to few-shot learning to get various metrics: CoSP-0, CoSP-1, SSV (the sum of averaged shapley value, i.e. $\sum_{i=1}^{n} \bar{\phi}(i)$), NoE(number of expressions, i.e. *n*). The correlations of these metrics and model inference accuracy across diverse datasets and models in 1, 2, 4-shot scenarios are shown in Tab. 1, and Tab. 2 record the correlations averaged among different models.

Observed from Tab. 1, CoSP-0 is the best interpretation metric for all models, and the interpretation validity of CoSP-0/CoSP-1 is much better than the other metrics. According to Tab. 2, CoSP-0 serves as the best interpretation metric for GSM8K, MultiArith, and SVAMP, while CoSP-1 for AQUA. For MathQA, CoSP-0 serves as the best interpretation metric in 1 or 2-shot learning, while CoSP-1 the best in 4-shot learning.

4.3 Explanation Reliability: Large-Scale Testing Assessment of CoSP Explanations

To further assess the reliability of our CoSP explanations, we conducted comprehensive validation experiments using the entire test set of the GSM8K dataset with both the LLaMA 2 and LLaMA 3 models. This focused approach ensures generality while maintaining computational feasibility. We selected four demonstrations for each model where the CoSP-0 scores for LLaMA 2 is 173, 121, 280, 235, while for LLaMA 3 is 264, 220, 344, 334.

The experimental outcomes consistently demonstrated a strong positive correlation between CoSP-0 scores and model accuracy for both LLaMA 2 and LLaMA 3. Specifically, for LLaMA 2, the demon-



set across different models and demos: (a) LLaMA2-demo1, (b) LLaMA3-demo1, (c) LLaMA2-demo2, and (d) LLaMA3-demo2. The observed results indicate that the accuracy curve for low CoSP expressions encompasses that for high CoSP expressions in almost all scenarios, highlighting that alterations on low CoSP expressions yield overall better performance outcomes compared to alterations on high CoSP expressions.

stration with a CoSP-0 score of 280 achieved the highest accuracy, followed by demonstrations with scores of 235, 173, and 121, in descending order of performance. Similarly, for LLaMA 3, the demonstration with a CoSP-0 score of 344 yielded the highest accuracy, followed by those with scores of 334, 264, and 220. This consistent pattern across both models indicates that demonstrations with higher CoSP-0 scores significantly enhance the reasoning capabilities of the models, while those with lower scores contribute less effectively.

505

510

511

512

513

514

515

516

517

518

519

521

522

525

527

530

531

534

538

4.4 Analytical Extensibility: Discovery of Novel Insights in CoT

Building upon our previous findings that high CoSP expressions contribute maximally, while low ones contributes minimally to model reasoning, we sought to uncover novel insights into the dynamics of CoT reasoning processes. Specifically, we applied four distinct altering to the expression with highest and lowest CoSP to assess their impact on model performance. 1) Remove the expression. 2) Replaced the expressions with non-informative placeholders, i.e. '...'. 3) Introduced calculation errors, for example, converting from 2 + 3 = 5 to 2+3 = 6'. 4) Introduced process errors, for example, converting from 2 + 3 = 5 to 4 + 7 = 11. And we selected two demononstrations and conducted these experiments on GSM8K datasets, with both the LLaMA 2 and LLaMA 3 models. The original demonstration is presented in Appendix B, where different expressions of CoSP in different colors.

Figures 3 depict the effect of these alterations on the accuracy of the test set for low and high CoSP expressions across different demonstrations and models. It was consistently observed across almost all experiments that the performance curves for low CoSP expressions encapsulated those for high CoSP expressions.

The results suggest that modifications to low CoSP expressions lead to better performance outcomes compared to modifications to high CoSP expressions. This finding further corroborates our initial hypothesis: low CoSP expressions exert minimal influence on model reasoning, whereas high ones significantly contribute.

Additionally, our experimental findings reveal several intriguing phenomena. Notably, the removal of certain expressions, the substitution of expressions with non-informative filler tokens (such as '...'), and the introduction of errors in either the result or process of expressions do not necessarily lead to significant degradation in model performance. This outcome resonates with prior studies(Pfau et al., 2024; Wang et al., 2023).

5 Conclusion

In this paper, we propose SalaMAnder, a novel framework for understanding and optimizing Chain-of-Thought (CoT) reasoning in large language models (LLMs). By introducing a theoretically grounded methodology based on Shapley value attribution and developing the CoSP (Cardinality of Shapley Positives) metric, we have established a mathematically rigorous approach to quantifying component-level contributions in CoT reasoning. Extensive validation across various LLM models and mathematical benchmarks demonstrates that the CoSP metric within our Sala-MAnder framework strongly and monotonically correlates with model performance. This correlation not only theoretically explains the empirical success of existing CoT but also provides rigorous guidelines for optimizing prompt construction. Furthermore, it can be utilized to discover novel insights resonating with prior studies

575

539

540

541

542

543

Limitation

References

and do case studies.

The limitation of this paper is that we do not have

enough time to conduct all complete experiments

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

2357–2367, Minneapolis, Minnesota. Association for

Jasmijn Bastings and Katja Filippova. 2020. The ele-

phant in the interpretability room: Why use attention

as explanation when we have saliency methods? In

Proceedings of the Third BlackboxNLP Workshop

on Analyzing and Interpreting Neural Networks for

NLP, pages 149-155, Online. Association for Com-

Xinyun Chen, Ryan A. Chi, Xuezhi Wang, and Denny

Zhou. 2024. Premise Order Matters in Reason-

ing with Large Language Models. arXiv e-prints,

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,

Abhinav Pandey, Abhishek Kadian, Ahmad Al-

Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-

ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh

Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-

tra, Archie Sravankumar, Artem Korenev, Arthur

Hinsvark, Arun Rao, Aston Zhang, Aurelien Ro-

driguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern,

Charlotte Caucheteux, Chaya Nayak, Chloe Bi,

Chris Marra, Chris McConnell, Christian Keller,

Christophe Touret, Chunyang Wu, Corinne Wong,

Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-

lonsius, Daniel Song, Danielle Pintz, Danny Livshits,

Danny Wyatt, David Esiobu, Dhruv Choudhary,

Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,

Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,

Elina Lobanova, Emily Dinan, Eric Michael Smith,

Filip Radenovic, Francisco Guzmán, Frank Zhang,

Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mi-

alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,

Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan

Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-

han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,

Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,

Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,

Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang,

Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,

Computational Linguistics.

putational Linguistics.

arXiv:2402.08939.

577 578

- 583 584 585
- 586
- 589
- 591

- 598
- 600

602

604 606 607

610 611 612

613

614

619

621

627

622 623

624

631

Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, and Tobias Speckbacher. 2024. The Llama 3 Herd of Models. arXiv e-prints, arXiv:2407.21783.

632

633

634

635

636

637

638

639

640

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

- Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao, Wenyue Hua, Yanda Meng, Yongfeng Zhang, and Mengnan Du. 2024. The Impact of Reasoning Step Length on Large Language Models. arXiv e-prints, arXiv:2401.04925.
- Jiachun Li, Pengfei Cao, Chenhao Wang, Zhuoran Jin, Yubo Chen, Daojian Zeng, Kang Liu, and Jun Zhao. 2024. Focus on Your Question! Interpreting and Mitigating Toxic CoT Problems in Commonsense Reasoning. arXiv e-prints, arXiv:2402.18344.
- Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. 2023. Symbolic Chain-of-Thought Distillation: Small Models Can Also "Think" Step-by-Step. arXiv e-prints, arXiv:2306.14050.
- 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.
- Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su. 2023.

- 691Deductive Verification of Chain-of-Thought Reason-692ing. arXiv e-prints, arXiv:2306.03872.
 - Andreas Madsen, Nicholas Meade, Vaibhav Adlakha, and Siva Reddy. 2022. Evaluating the faithfulness of importance measures in NLP by recursively masking allegedly important tokens and retraining. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1731–1751, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
 - Hoang H. Nguyen, Ye Liu, Chenwei Zhang, Tao Zhang, and Philip S. Yu. 2023. CoF-CoT: Enhancing Large Language Models with Coarse-to-Fine Chain-of-Thought Prompting for Multi-domain NLU Tasks. arXiv e-prints, arXiv:2310.14623.

703

706

708

711

712

713

714

715

718

719

720

721

723

727

729

731

732

733

734

735

736

737

739

740

741

742

743

744

745

- OpenAI. 2024. Openai o1 system card. [Online]. https://cdn.openai.com/ o1-system-card-20241205.pdf.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.
 - Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are NLP models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics.
 - Jacob Pfau, William Merrill, and Samuel R. Bowman. 2024. Let's Think Dot by Dot: Hidden Computation in Transformer Language Models. *arXiv e-prints*, arXiv:2404.15758.
 - Sofia Serrano and Noah A. Smith. 2019. Is attention interpretable? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2931–2951, Florence, Italy. Association for Computational Linguistics.
 - Lloyd S Shapley. 1953. A value for n-person games. Contributions to the Theory of Games, 2:307–317.
- Shoaib Ahmed Siddiqui, Radhika Gaonkar, Boris Köpf, David Krueger, Andrew Paverd, Ahmed Salem, Shruti Tople, Lukas Wutschitz, Menglin Xia, and Santiago Zanella Béguelin. 2024. Permissive information-flow analysis for large language models. *ArXiv*, abs/2410.03055.
- Qwen Team. 2024. Qwen2.5 technical report. *arXiv* preprint arXiv:2412.15115.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. arXiv e-prints, arXiv:2307.09288.

746

747

749

750

753

754

755

756

757

758

759

760

761

764

766

767

768

769

770

771

774

775

778

781

782

783

784

785

787

788

790

791

792

793

794

795

796

797

798

799

800

801

802

803

- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2023. Towards understanding chain-of-thought prompting: An empirical study of what matters. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2717–2739, Toronto, Canada. Association for Computational Linguistics.
- Lei Wang, Dongxiang Zhang, Lianli Gao, Jingkuan Song, Long Guo, and Heng Tao Shen. 2018. Mathdqn: solving arithmetic word problems via deep reinforcement learning. In Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence, AAAI'18/IAAI'18/EAAI'18. AAAI Press.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-Consistency Improves Chain of Thought Reasoning in Language Models. *arXiv e-prints*, arXiv:2203.11171.
- Robert James Weber. 1988. *Probabilistic values for* games, page 101–120. Cambridge University Press.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic Chain of Thought Prompting in Large Language Models. *arXiv e-prints*, arXiv:2210.03493.

- 805 Chuanyang Zheng, Zhengying Liu, Enze Xie, Zhenguo
 806 Li, and Yu Li. 2023. Progressive-Hint Prompting Improves Reasoning in Large Language Models. *arXiv*808 *e-prints*, arXiv:2304.09797.
- Bonny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, and Ed Chi.
 2022. Least-to-Most Prompting Enables Complex Reasoning in Large Language Models. arXiv eprints, arXiv:2205.10625.

817

818

821

822

825

826

827

830

A The proof of Theorems

- 816 We have three assumptions necessary for the proof:
 - 1. The positive contribution of any expression has a significant lower bound:

819
$$\exists \delta_+ > 0, \ s.t.$$

820 $\mu_i > \delta_+ \cdot \mathbb{I}(\mu_i > 0)$

2. The non-positive contribution of any expression has a lower bound:

823
$$\exists \delta_{-} < 0, \ s.t.$$

824
$$\mu_{i} > \delta_{-} \mathbb{I}(\mu_{i} \leq 0) = \delta_{-} \cdot (1 - \mathbb{I}(\mu_{i} > 0))$$

 The contributions of different expressions are mutually independent when applied to different problems:

828
$$\operatorname{Cov}(\phi^{(k)}(i), \phi^{(l)}(j)) = 0$$
829
$$(\forall i \neq j, 1 \leq k, l \leq m, k \neq l)$$

Here is the proof of Theo. 1:

Proof 1 As illustrated in Sec. 3.3:

832
$$\phi^{(k)}(i) \sim \mathcal{N}(\mu_i, \sigma_i^2)$$
833
$$\bar{\phi}(i) = \frac{1}{m} \sum_{k=1}^m \phi^{(k)}(i) \xrightarrow{m \to \infty} \mu_i$$

To simplify the expression, we define a positive contribution indicator $X_i = \mathbb{I}(\bar{\phi}(i) > 0)$. Thus:

836
$$CoSP = \sum_{i=1}^{n} X_i - \lambda \sum_{i=1}^{n} (1 - X_i)$$
837
$$= (1 + \lambda) \sum_{i=1}^{n} X_i - n\lambda$$
(7)

And we define the model performance Perf by
summing the expected shapley value of all expressions:

841
$$Perf = \sum_{i=1}^{n} \mathbb{E}[\phi(i)] = \sum_{i=1}^{n} \mu_i$$
(8)

Thus we can further derive the expression of Perf: 842

$$Perf = \sum_{i \in S_{+}} \mu_{i} + \sum_{i \notin S_{+}} \mu_{i} > \sum_{i \in S_{+}} \delta_{+} + \sum_{i \notin S_{+}} \delta_{-}$$
843

$$=\sum_{i=1}^{n}\delta_{+}\mathbb{I}(\mu_{i}>0)+\sum_{i=1}^{n}\delta_{-}\mathbb{I}(\mu_{i}\leqslant0)$$
844

$$=\sum_{i=1}^{n} \delta_{+} \mathbb{I}(\mu_{i} > 0) + \sum_{i=1}^{n} \delta_{-}(1 - \mathbb{I}(\mu_{i} > 0))$$
 845

$$= n\delta_{-} + \sum_{i=1}^{n} (\delta_{+} - \delta_{-}) \cdot \mathbb{I}(\mu_{i} > 0)$$
846

847

856

859

860

862

865

866

871

$$= n\delta_{-} + (\delta_{+} - \delta_{-}) \cdot \frac{CoSP + n\lambda}{1 + \lambda}$$
(9)

indicating a linear functional relationship between848a lower bound of model performance and CoSP.849And the coveriance between μ_i and X_i is:850

$$Cov(\mu_i, X_i) = \mathbb{E}[\mu_i X_i] - \mathbb{E}[\mu_i] \mathbb{E}[X_i]$$
⁸⁵

 $\begin{array}{ll} \mbox{where } \mu_i > \delta_+ X_i + \delta_- (1-X_i) \mbox{ based on the first} & \mbox{852} \\ \mbox{two assumptions.} & \mbox{853} \\ \mbox{We define a residual item } \epsilon_i > 0, \ s.t.: & \mbox{854} \end{array}$

$$\mu_i = \delta_+ X_i + \delta_- (1 - X_i) + \epsilon_i \tag{855}$$

Then

$$\mathbb{E}[\mu_i X_i] = \delta_+ \mathbb{E}[X_i^2] + \delta_- \mathbb{E}[(1 - X_i)X_i] + \mathbb{E}[\epsilon_i X_i]$$

$$= \delta_+ + \mathbb{E}[\epsilon_i X_i]$$
858

The second equation is because $X_i(1 - X_i) = 0$. And

$$\mathbb{E}[\mu_i] = \delta_+ \mathbb{E}[X_i] + \delta_- \mathbb{E}[1 - X_i] + \mathbb{E}[\epsilon_i]$$

Thus

$$Cov(\mu_i, X_i) = \delta_{+} \mathbb{E}[X_i^2] + \mathbb{E}[\epsilon_i X_i] - \delta_{+} \mathbb{E}^2[X_i] - \delta_{-} \mathbb{E}[X_i] \mathbb{E}[1 - X_i] + \mathbb{E}[X_i] \mathbb{E}[\epsilon_i]$$

Since $\mathbb{E}[X_i] = \mathbb{E}[X_i^2]$, and $\mathbb{E}[1 - X_i] = 1 - \mathbb{E}[X_i]$, then

$$\mathbb{E}[X_i]\mathbb{E}[1-X_i] = \mathbb{E}[X_i](1-\mathbb{E}[X_i])$$

$$= \mathbb{E}[X_i] - \mathbb{E}^2[X_i]$$
868

$$= \mathbb{E}[X_i^2] - \mathbb{E}^2[X_i]$$
869

$$= Var(X_i)$$
870

Then

$$Cov(\mu_i, X_i) = (\delta_+ - \delta_-) Var(X_i) + Cov(\epsilon_i, X_i)$$
(10)

(10)

Based on the third assumption, we have:

874
$$Cov(Perf, CoSP) = \sum_{i=1}^{n} \sum_{j=1}^{n} Cov(\mu_i, (1+\lambda)X_j - \lambda)$$

75
$$= \sum_{i=1}^{n} Cov(\mu_i, (1+\lambda)X_i - \lambda)$$

 $= (1 + \lambda) \sum_{i=1}^{n} Cov(\mu_i, X_i)$ 876

$$= (1+\lambda) \left[(\delta_{+} - \delta_{-}) \sum_{i=1}^{n} Var(X_{i}) + \sum_{i=1}^{n} Cov(\epsilon_{i}, X_{i}) \right]$$

Thus the expected value of CoSP with n + 1expressions is:

$$\mathbb{E}[CoSP_{n+1}] = (1+\lambda)\sum_{i=1}^{n+1} p_i - (n+1)\lambda$$
901

$$= \mathbb{E}[CoSP_n] + p_{n+1} - \lambda \quad (16)$$

Therefore, CoSP-0 increases monotonically with 903 the number of expressions n, while CoSP-1 de-904 creases monotonically with n. 905

878

879

884

886

887

888

890

891

895 896

898

8

And since the residual ϵ_i has little relevance with X_i , the sum of the covariance tends to 0. Thus

880
$$Cov(Perf, CoSP) = (1 + \lambda)(\delta_{+} - \delta_{-}) \sum_{i=1}^{n} Var(X_{i})$$
881
$$> 0$$
(11)

Specifically, we define CoSP-0 and CoSP-1, with 882 λ equals to 0 and 1, respectively. Then 883

$$Cov(Perf, CoSP-0) = (\delta_{+} - \delta_{-}) \sum_{i=1}^{n} Var(X_{i})$$
(12)

885
$$Cov(Perf, CoSP-I) = 2(\delta_{+} - \delta_{-}) \sum_{i=1}^{n} Var(X_{i})$$
(13)

Thus CoSP has a positive correlation with model performance.

Here is the proof of Theo. 2:

Proof 2 Since $X_i = \mathbb{I}(\overline{\phi}(i) > 0)$, then X_i follows a Bernoulli distribution:

$$p_i = P(X_i = 1) = \Phi(\frac{\mu_i}{\sigma_i}) \tag{14}$$

where $\Phi(\cdot)$ is the standard normal distribution cu-893 mulative function. 894

> Thus the expected value of CoSP with n expressions is:

897

$$\mathbb{E}[CoSP_n] = (1+\lambda)\sum_{i=1}^n \Phi(\frac{\mu_i}{\sigma_i}) - n\lambda$$
$$= (1+\lambda)\sum_{i=1}^n p_i - n\lambda \qquad (15)$$

899

900

908

909

910

911

912

913

B Selected Demonstrations

This section presents the selected demonstrations in Sec 4.4. Expressions with a light blue background have the lowest CoSP, those with an orange background have the highest CoSP, and the remaining expressions are shown with a light green background.

demo1

Question:

Sharon wants to get kitchen supplies. She admired Angela's kitchen supplies which consist of: 20 pots, 6 more than three times as many plates as the pots, and half as many cutlery as the plates. Sharon wants to buy: half as many pots as Angela, 20 less than three times as many plates as Angela, and twice as much cutlery as Angela. What is the total number of kitchen supplies Sharon wants to buy? Answer: Angela has 6+3*20=«6+3*20=66»66 plates. Angela has 1/2*66=«1/2*66=33»33 cutlery. Sharon wants to buy 1/2*20=«1/2*20=10»10 pots. Sharon wants to buy

 $3*66-20=x^3*66-20=178$ »178 plates. Sharon wants to buy $2*33=x^2*33=66$ »66 cutlery. Sharon wants to buy a total of $10+178+66=x^{10}+178+66=254$ »254 kitchen supplies.

demo2

Question:

Brittany, Alex, and Jamy all share 600 marbles divided between them in the ratio 3:5:7. If Brittany gives Alex half of her marbles, what's the total number of marbles that Alex has?

Answer:

The total ratio representing the number of marbles is 3+5+7 = (3+5+7) = 15the From the fraction representing ratio. the number of marbles that Brittany is 3/15, which is has equal

to	3/15*600 = «3/15*600=1	20»120						
marbles.Alex has								
5/15*60	00 = «5/15*600=200»200							
marbles	.If Brittany gives	half						
of her marbles to Alex, Alex receives								
<pre>1/2*120 = 60 marbles.After receiving</pre>								
60 marb	les from Brittany, Al	ex has						
<mark>200+60 = «200+60=260»260</mark> marbles.								

916