CARMA: Enhanced Compositionality in LLMs via Advanced Regularisation and Mutual Information Alignment

Anonymous ACL submission

Abstract

Large language models (LLMs) struggle with compositional generalisation, limiting their ability to systematically combine learned components to interpret novel inputs. While architectural modifications, fine-tuning, and data augmentation improve compositionality, they often have limited adaptability, face scalability constraints, or yield diminishing returns on real data. To address this, we propose CARMA, an intervention that enhances the stability and robustness of compositional rea-011 soning in LLMs while preserving fine-tuned 012 performance. CARMA employs mutual in-014 formation regularisation and layer-wise stability constraints to mitigate feature fragmentation, ensuring structured representations persist across and within layers. We evaluate CARMA on inverse dictionary modelling and 019 sentiment classification, measuring its impact on semantic consistency, performance stability, and robustness to lexical perturbations. Results show that CARMA reduces the variability introduced by fine-tuning, stabilises token representations, and improves compositional reasoning. While its effectiveness varies across architectures, CARMA's key strength lies in reinforcing learned structures rather than introducing new capabilities, making it a scalable auxiliary method. These findings suggest that integrating CARMA with fine-tuning can improve compositional generalisation while maintaining task-specific performance in LLMs.

1 Introduction

034

042

Compositional generalisation (CG) refers to the ability to systematically combine known expressions to generate novel ones following learned rules (Partee, 1984). This capability is essential for advancing language models (LMs) towards robust linguistic understanding beyond mere pattern matching (Ram et al., 2024).

Despite their strong performance across various natural language processing tasks, large language

models (LLMs) exhibit persistent weaknesses in compositional generalisation (Hupkes et al., 2020; Kim and Linzen, 2020a; Aljaafari et al., 2024). These limitations stem from multiple factors, including training objectives and model architectures. Standard autoregressive training methods, such as next-token prediction, prioritise statistical correlations in token sequences over structured semantic understanding (Yin et al., 2023a; Dziri et al., 2024). As a result, token representations often lack structured compositionality, leading to fragmented information processing within layers (horizontal misalignment) and across layers (vertical inconsistency). 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

Additionally, while self-attention mechanisms in Transformer models effectively capture local dependencies, they frequently fail to maintain coherent compositional representations across multiple layers (Murty et al., 2023). This misalignment impairs the model's ability to generalise compositionally, resulting in sensitivity to input order (Ismayilzada et al., 2024) and difficulties in handling complex syntactic and morphological structures (Aljaafari et al., 2024).

Several approaches have been proposed to address these limitations, including architectural modifications, enhanced encoding strategies, and targeted regularisation techniques (Ontanon et al., 2022; Murty et al., 2023; Csordás et al., 2021). However, these methods often struggle to balance compositional improvements with maintaining performance across diverse downstream tasks. Moreover, their effectiveness is typically confined to specific compositional structures or synthetic benchmarks. Developing a robust and adaptable solution that enables LLMs to achieve consistent CG across diverse tasks remains a major challenge.

This work introduces **CARMA**: enhanced Compositionality in LLMs via Advanced Regularisation and Mutual Information Alignment, illustrated in Figure 1. CARMA enhances CG by



Figure 1: This diagram depicts the computation of the loss and illustrates the integration of the Mutual Information (MI) loss (\mathcal{L}_{MI}) and the Stability Loss ($\mathcal{L}_{stability}$) into the final optimisation process. Tokens Tok_1 and Tok_2 form the positive set (H_{pos}) , while Tok_3, Tok_4, Tok_5 form the negative set (H_{neg}) . The \mathcal{L}_{MI} loss is computed vertically across layers (l to k), maximising the similarity of tokens in H_{pos} while contrasting them with tokens in $H_{\text{neg.}}$ The $\mathcal{L}_{\text{stability}}$ loss is computed *horizontally* between consecutive layers, ensuring consistency in hidden state representations. Both auxiliary losses are combined with the task loss (\mathcal{L}_{task}) to form the total loss (\mathcal{L}_{total}). This integration improves token representations and enhances the model's overall optimisation.

addressing training challenges that hinder structured compositionality in LLMs. By balancing layer-specific updates and reinforcing token-level dependencies, CARMA provides a scalable and adaptable solution that improves CG without sacrificing downstream task performance. То evaluate CARMA's effectiveness, we investigate the following research questions:

> • RQ1: How does regulating mutual information across layers influence compositionality in LLMs? How does it affect sensitivity to input and internal perturbations?

094

100

101

102

103

106

107

108

111

• RO2: To what extent does layer-specific regularisation improve compositional generalisation across semantic and sentiment analysis tasks, assessing CARMA's adaptability across domains?

The key contributions of this work are as follows:

- A novel regularisation method that enhances compositional generalisation without requiring architectural modifications. CARMA leverages mutual information alignment to preserve token dependencies across layers and employs layer-wise stability constraints to reduce representational inconsistencies.
- A systematic evaluation of CARMA across 109 compositionally demanding tasks, demonstrat-110 ing its ability to reinforce systematicity and

substitutivity, particularly in models where fine-tuning alone is insufficient.

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

• A theoretical and empirical analysis of how token dependencies degrade across layers in standard LLMs, revealing that CG limitations are not solely dependent on model size but rather on representational instability. CARMA mitigates this by ensuring consistent information flow, showing that non-intrusive regularisation strategies can significantly improve CG.

The remainder of this paper is structured as follows: Section 2 reviews compositionality in LLMs and associated challenges. Section 3 introduces the CARMA method. Section 4 describes the experimental setup. Section 5 presents empirical findings. Section 6 discusses related work. Section 7 offers insights and future research directions. Supporting datasets and software are available at a public repository.¹

2 **Compositionality in LLMs**

Compositional generalisation (CG) in linguistics encompasses five key principles: systematicity, productivity, substitutivity, localism, and overgeneralisation (Dankers et al., 2022a). These principles have been explored in LLMs for various applications, including compositional instruction (Yang

¹Anonymised for review.

228

229

230

231

232

233

234

235

236

238

et al., 2024b), semantic parsing (Li et al., 2023), machine translation (Li et al., 2021), and multi-step inference (Zhang et al., 2024). Empirical studies reveal that standard Transformer-based LLMs exhibit limited CG, even for relatively simple compositional tasks. For instance, models frequently struggle to assemble tokens into words or construct morphemes into coherent structures (Aljaafari et al., 2024; Ismayilzada et al., 2024). These limitations are linked to architectural constraints, training objectives, and tokenisation practices, which fragment information and increase sensitivity to input order and contextual noise (Murty et al., 2023).

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

155

156

157

158

159

162

163

165

167

168

169

170

171

173

174

175

176

177

178

179

180

182

186

187

190

Training Objectives and Information Fragmentation. Standard training objectives for LLMs typically optimise for next-token prediction, which prioritises surface-level correlations over deeper semantic integration (Dziri et al., 2024). While this approach is effective for data already seen, it often impedes CG by reducing mutual information between dependent tokens, thereby limiting the model's ability to form coherent compositional representations (Aljaafari et al., 2024).

Architectural Mechanisms and Compositional Consistency. Beyond training objectives, architectural mechanisms such as dropout and selfattention contribute to the dispersion of information across the model. This fragmentation increases sensitivity to input order and context, often resulting in errors that undermine compositional consistency (Sajjadi et al., 2016; Cai et al., 2021)—the model's ability to maintain produce consistent outputs when processing variations of semantically equivalent inputs through transformations like word substitution or paraphrasing.

These challenges impact both high-complexity reasoning tasks and simpler operations that demand consistent morphological and syntactic processing (Ismayilzada et al., 2024).

Existing Approaches to Enhance CG in LLMs. To address CG limitations, research has explored architectural adjustments, regularisation techniques, and task-specific strategies. For instance, (Ontanon et al., 2022) demonstrated that combining relative positional encoding with embeddings enhances CG, particularly in algorithmic tasks. Their findings suggest that weight sharing and copy decoders help retain input structures, thus improving CG accuracy. Other architectural modifications, such as Pushdown Layers (Murty et al., 2023) and GroCoT (Sikarwar et al., 2022), incorporate mechanisms for tracking syntactic depth and spatial relations,

which enable recursive processing of compositional structures.

Models like RegularGPT (Chi et al., 2023) introduce adaptive depth and memory mechanisms to facilitate CG by constructing complex structures from simpler components. Studies by (Csordás et al., 2021) and (Petty et al., 2024) evaluate model depth, parameter configurations, and encoding methods, revealing that architectural choices and training setups-such as avoiding early stopping and prioritising accuracy over loss minimisation-are critical to enhancing CG. In neural machine translation (NMT), (Dankers et al., 2022b) reformulated CG evaluations, finding a positive correlation between data size and compositional performance, underscoring the importance of extensive, real-world benchmarks for capturing the complexities of linguistic compositionality.

Frameworks like CompMCTG and Meta-MCTG (Zhong et al., 2024) offer benchmarks for evaluating CG in multi-aspect text generation, suggesting that joint training and meta-learning approaches can improve fluency. However, significant performance drops persist in out-of-distribution tasks. Additionally, synthetic tasks reveal that recursive, step-by-step prompt formats support combinatorial generalisation, although training biases and sequence order constraints remain limiting factors (Ramesh et al., 2024).

Enhanced Compositionality via 3 **Advanced Regularisation and Mutual Information Alignment (CARMA)**

This section formalises compositionality, introduces the core principles of CARMA, and details its components. Figure 1 illustrates the CARMA method, highlighting its optimisation process and key components.

3.1 **Compositionality Formalisation**

Mathematical Foundations of Compositionality. CG (Section 2) can be formally defined through a compositional system where \mathcal{E} denotes a set of expressions (e.g., token sequences recognised by the model), and $\mathcal M$ represents a corresponding set of meanings. This relationship is formalised as a function:

$$f: \mathcal{E} \to \mathcal{M} \tag{1}$$

For any complex expression $e \in \mathcal{E}$, composed of 237 constituent elements e_1, \ldots, e_n according to a syn-

241

242 243

244 245

- 247

248 249

251

255

258

261

263

265

266

267

270

271

273

274

3.2 CARMA Formalisation

maintaining interpretability.

CARMA operates over a range of target layers, from l to \mathcal{K} ($0 < l \leq \mathcal{K} \leq L$, where L is the total number of layers), and consists of two core components: Mutual Information and Layer-Wise Stability Regularisation.

tactic rule r, the function f satisfies:

to the syntactic rule r.

bound or deviation α :

 $f(r(e_1, \ldots, e_n)) = g_r(f(e_1), \ldots, f(e_n)),$

where q_r is the semantic operation that corresponds

Compositional Generalisation in LLMs. Effec-

tive CG in LLMs requires generating structured

compositions that preserve semantic consistency.

Given a novel expression e_{novel} similar to a known

expression e_{known} within a threshold β , their seman-

tic functions must remain within an interpretable

This formulation captures systematicity (struc-

tured combinations), substitutivity (preservation under transformations), and resistance to over-

generalisation (bounded semantic deviation) while

(2)

(3)

Mutual Information (MI) Regularisation Across Layers. CARMA preserves essential dependencies and maintains structural coherence by maximising MI between hidden states of related tokens. The MI between hidden states h_i^k and h_i^k at layer k, representing two related tokens i and j, is defined as:

$$I(h_{i}^{k}; h_{j}^{k}) = \mathbb{E}_{P(h_{i}^{k}, h_{j}^{k})} \left[\log \frac{P(h_{i}^{k}, h_{j}^{k})}{P(h_{i}^{k})P(h_{j}^{k})} \right]$$
(4)

Since exact computation is intractable, MI is approximated using the InfoNCE loss (Oord et al., 2018), encouraging token-level dependencies across the same layers:

$$\mathcal{L}_{\mathrm{MI}} = -\frac{1}{N} \sum_{k=l}^{\mathcal{K}} \sum_{i=1}^{Q} \left(\log \sum_{\substack{h_j \in \mathcal{H}^k \\ j \neq i}} \exp\left(\frac{f(h_i^k, h_j^k)}{\tau}\right) - \log\left(\sum_{\substack{h_j \in \mathcal{H}^k \\ j \neq i}} \exp\left(\frac{f(h_i^k, h_j^k)}{\tau}\right) + \sum_{\substack{h_m \in \mathcal{N}^k}} \exp\left(\frac{f(h_i^k, h_m)}{\tau}\right) \right) \right),$$
(5)

where $f(h_i^k, h_j^k)$ is a similarity function quantifying the relationship between hidden states at layer k, \mathcal{H}^k denotes the set of positive examples related to h_i^k , \mathcal{N}^k is the set of negative examples unrelated to h_i^k at layer k, τ is the temperature parameter, and N is the total number of target layers from lto \mathcal{K} , with Q representing the number of tokens or samples used per layer. Further details on MI approximation are provided in Appendix D.

275

276

277

278

279

281

284

286

289

290

291

292

293

294

295

296

299

300

301

302

304

305

306

307

308

309

310

311

312

Layer-Wise Stability Regularisation. This component enforces smooth transitions across layers, reducing abrupt changes that could disrupt compositional structures. For a layer k, the Layer-Wise $d(e_{\text{novel}}, e_{\text{known}}) \leq \beta \Rightarrow d(f(e_{\text{novel}}), f(e_{\text{known}})) \leq \alpha$. Stability Loss is defined as:

$$\mathcal{L}_{\text{Stability}} = \sum_{k=l}^{\mathcal{K}} \mathbb{E} \left[\frac{\left| f^{(k+1)}(X) - f^{(k)}(X) \right|^2}{\mathbb{E} \left[\left| f^{(k)}(X) \right|^2 \right] + \mathbb{E} \left[\left| f^{(k+1)}(X) \right|^2 \right] + \epsilon} \right],$$
(6)

where $f^{(k)}(X)$ denotes the activation output at layer k, and ϵ is a small positive constant to ensure numerical stability (e.g., $\epsilon = 10^{-8}$). Minimising this loss preserves compositional integrity across the specified layers by encouraging smooth and consistent transitions between them, thereby enabling more stable information flow and aggregation within this range.

CARMA Loss. CARMA integrates \mathcal{L}_{MI} and $\mathcal{L}_{Stability}$ into its total loss as:

$$\mathcal{L}_{\text{CARMA}} = \gamma \mathcal{L}_{\text{MI}} + \eta \mathcal{L}_{\text{Stability}}, \qquad (7)$$

where γ and η are hyperparameters in [0, 2] that control the relative contribution of each component. The final optimisation objective balances task-specific performance with CARMA's regularisation as:

$$\mathcal{L}_{\text{total}} = (1 - \lambda) \cdot \mathcal{L}_{\text{task}} + \lambda \cdot \mathcal{L}_{\text{CARMA}}, \quad (8)$$

where \mathcal{L}_{task} represents the task-specific loss, \mathcal{L}_{CARMA} is the regularisation loss, and $\lambda \in [0, 2]$ controls the trade-off between task accuracy and compositional robustness.

Experimental Setup 4

Downstream Tasks & datasets 4.1

CARMA is evaluated across two tasks that assess 313 different aspects of compositional generalisation: 314 Inverse Dictionary Modelling for word-level com-315 position and Sentiment Classification for phrase-316 level structure. These tasks measure systematicity, 317

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

368

- substitutivity, over-generalisation, and robustnessto perturbations.
- **Inverse Dictionary Modelling (IDM)** evaluates a model's ability to generate terms from definitions, focusing on substitutivity in semantic composition. WordNet (Miller, 1994) is used as the 323 training dataset, with an 80-10-10 train-validation-324 test split. Models are prompted with a definition and tasked with generating the corresponding term (e.g., "The star around which the Earth orbits is 327 called" \rightarrow "Sun"). Performance is assessed using Exact Match Accuracy, which measures whether 329 the generated term precisely matches the expected 330 output. By mapping definitions to terms, this task provides a robust assessment of a model's ability 332 to perform compositional substitution.
- Sentiment Classification (SC) assesses the model's ability to infer sentiment from phrases and sentences, particularly focusing on sentiment shifts 336 and over-generalisation. The Stanford Sentiment Treebank (SST) (Socher et al., 2013) is used with its original dataset splits. Models predict sentiment labels given textual inputs (e.g., "A brilliant perfor-340 mance sentiment is" \rightarrow "positive"). Performance 341 is evaluated using Exact Match Accuracy. This 342 task examines how sentiment composition is preserved across different levels of linguistic structure. Task formalisation, dataset details, and task selection rationale are in Appendices A, B.1, and B.2, respectively.

4.2 Model Configurations and Baselines

351

352

365

367

Experiments are conducted across three model configurations: baseline models, models with taskspecific fine-tuning, and models with fine-tuning plus CARMA regularisation. Models use 500 warm-up steps and a 0.006 learning rate. We test GPT-2 (S/L) (Radford et al., 2019), Gemma-2B (Team et al., 2024), Llama (1B/3B) (Dubey et al., 2024), and Qwen (0.5B/3B) (Yang et al., 2024a), representing diverse architectures and capacities. CARMA regularisation is generally applied at approximately one-third of the model's depth, though specific layer positions vary. Details on fine-tuning methodologies, model specifications, and CARMA hyperparameter selection are provided in Appendix B.3.

4.3 Interventions for Compositional Robustness and Performance Stability

Two interventions are used to evaluate the robustness of compositional structures and the stability of learned representations: Constituent-aware pooling and synonym replacement. These interventions assess hierarchical dependencies and semantic consistency under controlled perturbations.

Constituent-Aware Pooling (CAP) (Aljaafari et al., 2024) groups token-level representations into higher-level semantic units (e.g., words, syntactic constituents) to assess hierarchical dependencies and how compositional structures are maintained across layers. In this paper, the token-to-word CAP is utilised. Model robustness is measured by monitoring performance metrics before and after applying CAP. Full methodology and formalisation are provided in Appendix C.1.

Synonym Replacement evaluates semantic consistency by substituting 25% and 40% of prompt words with synonyms within an interpretable bound (α). Experiments were repeated at least five times with different seeds for robustness and performance stability assessment; further details are in Appendix C.2.

4.4 Experimental setup

Experiments were conducted using NVIDIA RTX A6000 and A100 GPUs. The method was developed in Python (v3.10.15) with Transformers (v4.44.2) (Wolf et al., 2020), PyTorch (v2.4.1) (Paszke et al., 2019), and Transformer-lens (v2.8.1) (Nanda and Bloom, 2022). Preprocessing tasks, including tokenisation and tagging, used NLTK (v3.9.1) (Bird et al., 2009), spaCy (v3.7.2) (Honnibal et al., 2020), and TextBlob (v0.18.0) (Loria et al.), with Scikit-learn (v1.5.1) (Pedregosa et al., 2011) for evaluation.

5 Results and discussion

The method is evaluated across three aspects: (1) its impact on model robustness against compositionalbased perturbations, (2) its impact on model performance stability, and (3) its impact on model overall performance. See Appendix B.4 for a detailed breakdown of the evaluation metrics used for each aspect.

5.1 Constituent-Aware Pooling (CAP) Intervention

Figures 2(a) and 2(b) show the impact of CAP on IDM and SC tasks, comparing original, fine-tuned (FT) and CARMA models.² Model performance is

²Throughout this paper, models incorporating CARMA with FT are referred to as CARMA models.



Figure 2: Layer-wise performance comparison under CAP intervention, with performance averaged over three protocols (Mean CAP, Max CAP, Sum CAP) for Original, Fine-Tuned (FT), and CARMA (FT + CARMA) models. Layer numbers are normalised to their relative positions within each model to enable cross-architecture comparison. The IDM task (left) highlights CARMA's improvements in systematicity and stability, particularly in the early and middle layers. The SC task (right) demonstrates CARMA's ability to enhance robustness, though convergence with FT occurs in deeper layers.

averaged across three CAP protocols (Mean, Max, and Sum), with per-protocol results provided in Appendix E. The analysis examines how well models preserve compositionality under hierarchical pooling.

414

415

416

417

418

CARMA's effectiveness is influenced by model 419 size, tokenisation strategy, and task complexity. 420 In IDM tasks, CARMA models have consider-421 able gains when applying CAP at the earliest lay-422 ers (1% of model depth), particularly in models 423 with fine-grained tokenisation: Llama-1B (+3.61%)424 and Gemma-2B (+16.89%). GPT2-L, despite its 425 reliance on subword tokenisation, benefits from 426 CARMA over FT (+3.67%). However, Llama-3B 427 and Qwen-3B minimal improvements (+1.0%) sug-428 gest a capacity ceiling where increased model size 429 does not yield proportional gains due to training 430 data limitations. The combination of smaller scale 431 and multilingual training particularly affects Qwen-432 0.5B, where limited model capacity coupled with 433 broad language coverage appears to constrain En-434 glish-specific compositional learning, resulting in 435 reduced CARMA benefits. In SC tasks, tokeni-436 sation effects vary with task complexity. When 437 intervening at 25% layer position, Gemma-2B 438 and Llama-1B show the strongest gains (+27.38%, 439 440 +10.59%), while Llama-3B exhibits a marginal difference between CARMA and FT ($\sim 1\%$) but still 441 outperforms the Original model (+37.68%). These 442 results suggest that fine-tuning alone is sufficient 443 for simpler tasks, whereas structured interventions 444

like CARMA are particularly beneficial for more complex, compositional reasoning tasks.

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

In a layer-wise analysis, the impact of CARMA varies significantly across network depths, revealing crucial insights about compositional learning in transformers. Early layers (0-25%) benefit the most from regularisation, as they establish foundational compositional representations by exhibiting a weak notion of compositionality. Middle layers (25-75%) reinforce these patterns, maintaining structured feature dependencies with moderate improvements. Deeper layers (75-100%) show minimal benefits as the model transitions from compositional learning to task-specialised representations. This pattern aligns with previous findings on layer-wise compositional evolution in Transformers, where earlier layers capture hierarchical structure, while deeper layers exhibit increased task specificity (Feucht et al., 2024). CARMA can thus be strategically applied to control these early representations, maintaining beneficial compositional structure while allowing natural task-specific adaptations in deeper layers.

These findings demonstrate CARMA's effectiveness, particularly for models with granular tokenisation under data constraints, mediated by model capacity and task demands. The method's dual role - enhancing early compositional learning while preserving deeper layer adaptations - enables targeted improvement in model robustness without disrupting task-specific processing.

Model	Ver.	Task	Int.	CS	CV
GPT2-L	CARMA	IDM	25%	56.31	0.0164
	FT	IDM	25%	56.95	0.0311
	Org	IDM	25%	51.10	0.1175
	CARMA	SC	25%	0.8858	0.0065
	FT	SC	25%	0.8804	0.0082
	CARMA	IDM	25%	56.70	0.023
	FT	IDM	25%	57.42	0.030
Commo 2D	Org	IDM	25%	49.47	0.031
Gemma-2B	CARMA	SC	25%	78.90	0.008
	FT	SC	25%	80.23	0.009
	Org	SC	25%	68.14	0.042
Llama-3B	CARMA	IDM	25%	62.86	0.015
	FT	IDM	25%	62.22	0.029
	Org	IDM	25%	52.47	0.035
	CARMA	SC	25%	84.83	0.0056
	FT	SC	25%	85.85	0.0065
	Org	SC	25%	35.21	0.0136

Table 1: Model performance (25% synonym intervention). Ver.: Version; Int.: Intervention rate; CS: ConsistSyn (%); CV: Coefficient of Variation. Best values in bold.

5.2 Synonyms Replacement Intervention

476

477

478

479

480

481

482

483

484

485

486

487

Synonym Replacement evaluates semantic consistency and robustness under lexical variations across multiple runs ($N \ge 5$) with different seeds. *ConsistSyn* measures output preservation after substitution, while the coefficient of variation (CV) quantifies performance stability, with lower values indicating higher stability. Performance is assessed at 25% and 40% word replacement rates to measure sensitivity to increasing perturbations. A sample of results is presented in Table 1, with full details in Appendix E.

Across models, CARMA achieves a distinctive 488 performance profile, matching or exceeding FT 489 ConsistSyn while consistently demonstrating supe-490 rior stability through lower CV values. At 25% in-491 tervention, Gemma-2B CARMA achieves 56.70% 492 ConsistSyn with a CV of 0.0225, compared to 493 494 FT's 57.42% with higher variance (CV: 0.0307). Llama-3B CARMA outperforms FT in both Con-495 sistSyn (62.86% vs. 62.22%) and stability (CV: 496 0.0148 vs. 0.0292) for IDM. Qwen-3B follows a 497 similar trend but with smaller relative gains, im-498 proving stability (CV: 0.0225 vs. 0.0279) while 499 maintaining a marginal ConsistSyn advantage over 500 FT (62.00% vs. 61.79%). However, as interven-501 tion complexity increases to 40%, the performance gap widens; for example, Gemma-2B FT main-503 tains higher ConsistSyn (44.98%) than CARMA (42.36%), though CARMA remains more stable (CV: 0.0174 vs. 0.0249). This behaviour implies 507 that the advantage of CARMA lies in its lower variance and reinforcement of compositional con-508 sistency. Thus, it maintains compositional understanding without sacrificing performance, whereas 510 FT produces a performance-driven approach. 511

The tokenisation method significantly affects CARMA's impact. Models with more structured tokenisation show stronger stability improvements, but gains vary based on vocabulary design and language coverage. Llama and GPT2-L generally benefit more than Qwen, even with similar sizes, likely due to their smaller multilingual coverage, which results in a more compact and consistent token distribution. Qwen, with a larger vocabulary (151K tokens) supporting broader multilingual processing, introduces redundancy that dampens CARMA's relative stability advantage. Gemma-2B, optimised for a single dominant language with a large vocabulary size, shows the highest overall gains, reinforcing that a structured tokenisation approach focused on a limited linguistic scope enhances CARMA's effectiveness.

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

Task complexity further differentiates CARMA's effect. CARMA's advantages align with its methodological design, particularly in tasks requiring explicit structural reinforcement. In IDM, where systematicity and substitutivity are critical, CARMA ensures structured mappings hold under perturbation, particularly in Gemma-2B (+14.6% over the original) and Llama-1B (+2692.5% over the original in SC). However, in SC, where compositionality is more distributed, larger models show lower differences between CARMA and FT, reinforcing that larger models encode sentiment shifts effectively without additional intervention.

These results strengthen the hypothesis that CARMA enhances model robustness across perturbations, particularly in structured learning tasks and models where fine-tuning alone does not fully capture compositional dependencies. While FT maintains an advantage in absolute accuracy, CARMA ensures greater consistency, making it critical for improving compositional alignment and mitigating instability in high-variance settings.

5.3 Impact of CARMA on Performance

Figures 3 and 4 show the performance of original, FT, and CARMA accuracies across tasks. <u>CARMA</u> demonstrates significant improvements over original models across tasks. For example, in IDM, <u>GPT2-L achieves 150% improvement, and Llama-3B shows an 89.6% increase, while in SC, Gemma-2B demonstrates 122.5% improvement over Original baselines.</u>

Task-specific patterns emerge when comparing models. For instance, in IDM, CARMA outperforms FT, with Llama-3B showing a +5% gain



Figure 3: Task performance in IDM across GPT-2 (S, L), Gemma-2B, Llama (1B, 3B), and Qwen (0.5B, 3B).



Figure 4: Task performance in SC across GPT-2 (S, L), Gemma-2B, Llama (1B, 3B) and Qwen (0.5B, 3B).

and GPT2-L improving by 1.7%. In SC, CARMA maintains comparable performance to FT while enhancing robustness, suggesting it preserves learned features while strengthening compositional consistency.

CARMA enhances FT by improving representation stability and preventing feature drift, ensuring structured compositional consistency. Its benefits are most pronounced in larger models, where greater capacity supports robust representations while maintaining fine-tuned performance. This scalability highlights CARMA's effectiveness in regularising model representations and reinforcing compositional structure without disrupting learned task features, providing a reliable solution for improving compositional reasoning in LLMs.

6 Related work

563

564

569

571

572

577

580

581

582

Research on CG in LLMs has revealed both capabilities and limitations (Tull et al., 2024; Moisio et al., 2023; Sinha et al., 2024), though many studies lack mechanistic analysis or concrete suggestions for improvements.

585 Architectural modifications are a common ap-586 proach to tackle CG challenges. Recent proposals include pushdown layers for recursive attention (Murty et al., 2023), Layer-wise Representation Fusion for dynamic encoder weighting (Lin et al., 2023), and specialised semantic parsing methods (Shaw et al., 2021). While effective for specific tasks, these solutions face scalability challenges due to computational overhead, specialised annotation requirements, and architectural constraints. 587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

Regularisation methods provide alternative approaches through consistency regularisation (Yin et al., 2023b), data augmentation strategies (Ontanon et al., 2022), and attention stability mechanisms (Zhai et al., 2023). Studies show dataset complexity and example frequency variations improve compositional reasoning (Zhou et al., 2023). However, these methods face key limitations: tokenlevel approaches lack adaptability to complex structures, augmentation shows diminishing returns on real data, and stability mechanisms prioritise training stability over compositional generalisation.

Evaluation challenges persist in CG research. Standard benchmarks like SCAN (Lake and Baroni, 2017), PCFG (Hupkes et al., 2020), and COGS (Kim and Linzen, 2020b) rely heavily on synthetic data, limiting real-world applicability. Recent frameworks like CoGnition (Li et al., 2021) and CAP (Aljaafari et al., 2024) better align with natural language phenomena, but evaluation gaps remain. Current approaches often sacrifice generalisability for task-specific performance. *CARMA* addresses these limitations through a *task-agnostic*, *efficient solution* that enhances CG while maintaining robust cross-task performance.

7 Conclusion

This paper presents CARMA, a method for enhancing compositional generalisation in LLMs through mutual information regularisation and stability constraints. By addressing information fragmentation and layer-wise instability, CARMA improves performance stability and robustness under interventions, as demonstrated through IDM and SC tasks. The method offers a cost-effective solution applicable across model architectures with minimal modifications. Future work should explore extending CARMA to additional tasks that rely more on nuanced semantic features and multilingual settings to further evaluate its scalability and adaptability. Integrating CARMA into improved, targeted transformer architectures for CG could unlock new possibilities for enhancing compositionality.

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

686

637 Limitations

638The limitations of this paper can be summed up as639follows: First, our results are primarily reported640for the English language. Further analysis across641languages with diverse linguistic structures is left642as a confirmatory future work. Second, the datasets643(WordNet and SST) lack a more comprehensive644representativeness of broader linguistic phenomena.645Third, our focus is predominantly on decoder-based646Transformers. Finally, the employed Transformer647models may inherit potential biases ingrained from648their pre-training data.

Ethical statement

650This work aims to enhance language model ro-651bustness and compositional understanding through652CARMA. While improving model reliability is ben-653eficial, we acknowledge potential risks in enhanc-654ing language model capabilities. Our evaluation655focuses on controlled tasks (IDM and SC) with656comprehensive stability metrics to ensure responsi-657ble development and transparent reporting of model658behaviour under perturbations.

References

659

670

671

672

674

675

676

677

678

679

- Nura Aljaafari, Danilo S Carvalho, and André Freitas. 2024. Interpreting token compositionality in llms: A robustness analysis. *arXiv preprint arXiv:2410.12924*.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. Natural language processing with Python: analyzing text with the natural language toolkit. " O'Reilly Media, Inc.".
- Xingyu Cai, Jiaji Huang, Yuchen Bian, and Kenneth Church. 2021. Isotropy in the contextual embedding space: Clusters and manifolds. In *International conference on learning representations*.
- Ta-Chung Chi, Ting-Han Fan, Alexander Rudnicky, and Peter Ramadge. 2023. Transformer working memory enables regular language reasoning and natural language length extrapolation. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pages 5972–5984, Singapore. Association for Computational Linguistics.
- Róbert Csordás, Kazuki Irie, and Juergen Schmidhuber.
 2021. The devil is in the detail: Simple tricks improve systematic generalization of transformers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 619–634, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

- Verna Dankers, Elia Bruni, and Dieuwke Hupkes. 2022a. The paradox of the compositionality of natural language: A neural machine translation case study. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4154–4175, Dublin, Ireland. Association for Computational Linguistics.
- Verna Dankers, Elia Bruni, and Dieuwke Hupkes. 2022b. The paradox of the compositionality of natural language: A neural machine translation case study. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4154–4175. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. 2024. Faith and fate: Limits of transformers on compositionality. *Advances in Neural Information Processing Systems*, 36.
- Christiane Fellbaum. 1998. Wordnet: An electronic lexical database. *MIT Press google schola*, 2:678–686.
- Sheridan Feucht, David Atkinson, Byron Wallace, and David Bau. 2024. Token erasure as a footprint of implicit vocabulary items in llms. In *The 2024 Conference on Empirical Methods in Natural Language Processing*.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spacy: Industrialstrength natural language processing in python.
- Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. 2020. Compositionality decomposed: How do neural networks generalise? (extended abstract). In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 5065–5069. International Joint Conferences on Artificial Intelligence Organization. Journal track.
- Mete Ismayilzada, Defne Circi, Jonne Sälevä, Hale Sirin, Abdullatif Köksal, Bhuwan Dhingra, Antoine Bosselut, Lonneke van der Plas, and Duygu Ataman. 2024. Evaluating morphological compositional generalization in large language models. *arXiv preprint arXiv:2410.12656*.
- Najoung Kim and Tal Linzen. 2020a. Cogs: A compositional generalization challenge based on semantic interpretation. In *Proceedings of the 2020 conference on empirical methods in natural language processing (emnlp)*, pages 9087–9105.
- Najoung Kim and Tal Linzen. 2020b. COGS: A compositional generalization challenge based on semantic

Christopher Manning. 2023. Pushdown layers: Enence on Empirical Methods in Natural Language 799 Processing (EMNLP), pages 9087–9105, Online. Ascoding recursive structure in transformer language sociation for Computational Linguistics. models. In Proceedings of the 2023 Conference on 801 *Empirical Methods in Natural Language Processing*, 802 Nikita Kitaev, Steven Cao, and Dan Klein. 2019. Multipages 3233–3247, Singapore. Association for Com-803 lingual constituency parsing with self-attention and putational Linguistics. 804 pre-training. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Neel Nanda and Joseph Bloom. 2022. Transformerlens. 805 pages 3499–3505, Florence, Italy. Association for https://github.com/TransformerLensOrg/ 806 Computational Linguistics. TransformerLens. 807 Nikita Kitaev and Dan Klein. 2018. Constituency pars-Santiago Ontanon, Joshua Ainslie, Zachary Fisher, and 808 ing with a self-attentive encoder. In Proceedings Vaclav Cvicek. 2022. Making transformers solve 809 of the 56th Annual Meeting of the Association for compositional tasks. In Proceedings of the 60th An-810 Computational Linguistics (Volume 1: Long Papers), nual Meeting of the Association for Computational 811 pages 2676–2686, Melbourne, Australia. Association Linguistics (Volume 1: Long Papers), pages 3591– 812 for Computational Linguistics. 3607, Dublin, Ireland. Association for Computational 813 Linguistics. 814 Brenden M. Lake and Marco Baroni. 2017. Generalization without systematicity: On the compositional Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. 815 skills of sequence-to-sequence recurrent networks. Representation learning with contrastive predictive 816 In International Conference on Machine Learning. coding. arXiv preprint arXiv:1807.03748. 817 Yafu Li, Yongjing Yin, Yulong Chen, and Yue Zhang. Barbara H. Partee. 1984. Compositionality. In Fred 818 2021. On compositional generalization of neural ma-Landman and Frank Veltman, editors, Varieties of 819 chine translation. In Proceedings of the 59th Annual Formal Semantics, pages 281–312. Foris Publica-820 Meeting of the Association for Computational Lintions. 821 guistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Adam Paszke, Sam Gross, Francisco Massa, Adam 822 Papers), pages 4767-4780, Online. Association for Lerer, James Bradbury, Gregory Chanan, Trevor 823 Computational Linguistics. Killeen, Zeming Lin, Natalia Gimelshein, Luca 824 Antiga, et al. 2019. Pytorch: An imperative style, 825 Zhaoyi Li, Ying Wei, and Defu Lian. 2023. Learning high-performance deep learning library. Advances in 826 to substitute spans towards improving compositional neural information processing systems, 32. 827 generalization. In Proceedings of the 61st Annual Meeting of the Association for Computational Lin-F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, 828 guistics (Volume 1: Long Papers), pages 2791–2811, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, 829 Toronto, Canada. Association for Computational Lin-R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, 830 guistics. D. Cournapeau, M. Brucher, M. Perrot, and E. Duch-831 esnay. 2011. Scikit-learn: Machine learning in 832 Lei Lin, Shuangtao Li, Yafang Zheng, Biao Fu, Shan Python. Journal of Machine Learning Research, 833 Liu, Yidong Chen, and Xiaodong Shi. 2023. Learn-12:2825-2830. 834 ing to compose representations of different encoder layers towards improving compositional generaliza-Jackson Petty, Sjoerd Steenkiste, Ishita Dasgupta, Fei 835 tion. In Findings of the Association for Computa-Sha, Dan Garrette, and Tal Linzen. 2024. The impact 836 tional Linguistics: EMNLP 2023, pages 1599-1614, of depth on compositional generalization in trans-837 Singapore. Association for Computational Linguisformer language models. In Proceedings of the 2024 838 tics. Conference of the North American Chapter of the 839 Association for Computational Linguistics: Human 840 Steven Loria et al. textblob documentation. Release Language Technologies (Volume 1: Long Papers), 841 0.18.0. pages 7232-7245. 842 George A. Miller. 1994. WordNet: A lexical database Alec Radford, Jeffrey Wu, Rewon Child, David Luan, 843 for English. In Human Language Technology: Pro-Dario Amodei, Ilya Sutskever, et al. 2019. Language 844 ceedings of a Workshop held at Plainsboro, New models are unsupervised multitask learners. OpenAI 845 Jersey, March 8-11, 1994. blog, 1(8):9. 846 Anssi Moisio, Mathias Creutz, and Mikko Kurimo. 2023. Evaluating morphological generalisation in Parikshit Ram, Tim Klinger, and Alexander Gray. 2024. 847 machine translation by distribution-based composi-What makes models compositional? a theoretical 848 tionality assessment. In Proceedings of the 24th view. In Proceedings of the Thirty-Third Interna-849 Nordic Conference on Computational Linguistics tional Joint Conference on Artificial Intelligence 850 (IJCAI-24), pages 4824-4832. International Joint (NoDaLiDa), pages 738-751, Tórshavn, Faroe Is-851 lands. University of Tartu Library. Conferences on Artificial Intelligence Organization. 852

Shikhar Murty, Pratyusha Sharma, Jacob Andreas, and

798

742

743

745

746

747

748

752

753

757

758

759

762

770

771

773

774

777

778

779

781

782

784

790

791

793

794

797

interpretation. In Proceedings of the 2020 Confer-

Rahul Ramesh, Ekdeep Singh Lubana, Mikail Khona, Robert P. Dick, and Hidenori Tanaka. 2024. Compositional capabilities of autoregressive transformers: A study on synthetic, interpretable tasks. In *Forty-first International Conference on Machine Learning*.

853

857

858

859

861

862

864

865

866

867

870

871

874

876

878

879

881

882

885

889

895

900

901

902

903

904

905

906 907

- Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. 2016. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. *Advances in neural information processing systems*, 29.
- Peter Shaw, Ming-Wei Chang, Panupong Pasupat, and Kristina Toutanova. 2021. Compositional generalization and natural language variation: Can a semantic parsing approach handle both? In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 922–938, Online. Association for Computational Linguistics.
- Ankur Sikarwar, Arkil Patel, and Navin Goyal. 2022. When can transformers ground and compose: Insights from compositional generalization benchmarks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 648–669, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Sania Sinha, Tanawan Premsri, and Parisa Kordjamshidi. 2024. A survey on compositional learning of ai models: Theoretical and experimetnal practices. *arXiv preprint arXiv:2406.08787*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Sean Tull, Robin Lorenz, Stephen Clark, Ilyas Khan, and Bob Coecke. 2024. Towards compositional interpretability for xai. *arXiv preprint arXiv:2406.17583*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2019. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics. 908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024a. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.
- Haoran Yang, Hongyuan Lu, Wai Lam, and Deng Cai. 2024b. Exploring compositional generalization of large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 4: Student Research Workshop)*, pages 16–24.
- Yongjing Yin, Jiali Zeng, Yafu Li, Fandong Meng, Jie Zhou, and Yue Zhang. 2023a. Consistency regularization training for compositional generalization. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1294–1308.
- Yongjing Yin, Jiali Zeng, Yafu Li, Fandong Meng, Jie Zhou, and Yue Zhang. 2023b. Consistency regularization training for compositional generalization. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1294–1308, Toronto, Canada. Association for Computational Linguistics.
- Shuangfei Zhai, Tatiana Likhomanenko, Etai Littwin, Dan Busbridge, Jason Ramapuram, Yizhe Zhang, Jiatao Gu, and Joshua M Susskind. 2023. Stabilizing transformer training by preventing attention entropy collapse. In *International Conference on Machine Learning*, pages 40770–40803. PMLR.
- Min Zhang, Jianfeng He, Shuo Lei, Murong Yue, Linhan Wang, and Chang-Tien Lu. 2024. Can llm find the green circle? investigation and human-guided tool manipulation for compositional generalization. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 11996–12000. IEEE.
- Tianqi Zhong, Zhaoyi Li, Quan Wang, Linqi Song, Ying Wei, Defu Lian, and Zhendong Mao. 2024. Benchmarking and improving compositional generalization of multi-aspect controllable text generation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6486–6517. Association for Computational Linguistics.

966 967 968

969

971

972

973

974

975

976

978

979

980

982

983

985

991

995

997

999

1001

1002

1003

1004

1005

1006

1007 1008

1009

1010

1012

Xiang Zhou, Yichen Jiang, and Mohit Bansal. 2023. Data factors for better compositional generalization. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14549–14566, Singapore. Association for Computational Linguistics.

A Task Selection and Compositionality Considerations

To assess compositional generalisation and the benefits of CARMA, we targeted tasks that involve systematic meaning construction and sensitivity to structural modifications. To that end, we opted to employ Inverse Dictionary Modelling (IDM) and Sentiment Classification (SC) as proxies for different dimensions of compositionality, capturing both structured composition and hierarchical generalisation.

IDM requires models to generate a single-word representation from a natural language definition, mapping from the composition of input constituents (individual concept components) to a specific term. On the other hand, SC maps meaning to a sentiment label, aggregating local meaning elements into a global interpretation. While IDM focuses on explicit compositional mapping, SC evaluates distributed composition, where sentiment is shaped by multiple interacting components.

Both tasks assess several aspects of compositionality (Figure 5), namely systematicity (structured meaning formation), substitutivity (semantic preservation under transformation), and resistance to over-generalisation (ensuring bounded semantic deviation). Further, they evaluate robustness, testing whether models can maintain correctness and consistency under internal and input-lexical perturbations. IDM and SC provide a comprehensive test of compositional generalisation across structured and distributed representations.

B Detailed Experimental Configuration

B.1 Task Formalisation

This paper evaluates the effectiveness of CARMA in enhancing the compositional generalisation of large language models (LLMs) through two tasks. These tasks were selected based on their focus on input token structure and compositional semantics, utilising next-token prediction with single-token outputs. Formal definitions for each task are presented below.

1013 Inverse Definition Modelling (IDM). This task1014 requires the model to predict a definiendum D,



Figure 5: Illustration of compositional generalisation in Inverse Dictionary Modelling (IDM) and Sentiment Classification (SC). The figure highlights key compositional properties: systematicity ensures coherent meaning construction, substitutivity maintains meaning under lexical variations, robustness preserves intended outputs under perturbations, and over-generalisation leads to overly broad or semantically weak predictions (e.g., neuron misclassified as cell or positive reduced to neutral).

given its corresponding definition definition in natural language. Formally, the definition is represented as a sequence of tokens, definition = $\{tok_1, tok_2, ..., tok_n\}$, and the model seeks to produce D such that:

$$D = \arg\max_{t \in \mathcal{V}} P(d \mid \text{definition}), \qquad (9)$$

1015

1016

1017

1018

1019

1020

1021

1023

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

where \mathcal{V} denotes the model's vocabulary, and d represents a potential definiendum. Predictions are deemed correct only if they exactly match the target output.

Sentiment classification (SC). This task involves assigning a sentiment label to a given sentence containing sentiment cues and potential modifiers. The model processes the input sentence, represented as a sequence of tokens sentence = $\{tok_1, tok_2, ..., tok_n\}$, and produces an output label from a predefined set of sentiment classes \mathcal{A} (i.e., *positive*, *negative*, *neutral*). Formally, the task is defined as:

$$label = \arg \max_{\ell \in \mathcal{L}} P(\ell \mid sentence), \quad (10)$$

where $P(\ell \mid \text{sentence})$ is the probability of the sentiment label ℓ given the sentence. The model's 1035



Figure 6: IDM Performance Across Models Under CAP

1037performance is evaluated based on its ability to cor-1038rectly predict the sentiment, accounting for compo-1039sitional nuances such as modifiers and contrasts.

1041

1042

1043

1044

1045

1046

1047

1048

1051

1052

1053

1054

1055

1056

1057

1059

B.2 Datasets specification and pre-processing

For IDM, the training and test datasets were derived from WordNet (Fellbaum, 1998), a widely used lexical database of the English language. WordNet comprises over 117,000 synsets, each representing a distinct concept and annotated with semantic relationships such as hypernyms, synonyms, and definitions. To ensure consistency and improve data quality, standard preprocessing techniques were applied, including the removal of special characters, punctuation, extra spaces, and parenthesised content where necessary. The dataset focuses on general-purpose vocabulary rather than specialised domains or demographic groups. The dataset was initially split into an 80-20 ratio, with 80% allocated for training. The remaining 20% was further divided equally into validation and test sets.

> The SC dataset was derived from the Stanford Sentiment Treebank (SST) (Socher et al., 2013), a corpus of English movie reviews annotated for

analysis of the compositional effects of sentiment 1060 inference and was released under Apache License, 1061 Version 2.0. SST includes fine-grained sentiment 1062 labels at both the phrase and sentence levels, making it a standard benchmark for evaluating senti-1064 ment classification models. The original dataset 1065 splits provided by the authors were maintained to ensure consistency in training, validation, and test-1067 ing. For SST labels, sentiment scores were cate-1068 gorised as follows: values equal to or greater than 0.6 were classified as positive, scores between 0 1070 and 0.6 were considered neutral, and scores be-1071 low zero were assigned as negative. The final test dataset sizes for each task are presented in Table 2. 1073

Dataset	Train size	validation Size	Test Size
WordNet	9563	1154	1231%
SST	8544	1101	2210

Table 2: Train, validation, and test set sizes for WordNet and SST datasets used in this paper.

B.3 Model training and fine-tuning settings

1074

1076

 Table 3 summarises the key characteristics of the

 models evaluated in this study. All models were ob

tained from Hugging Face (Wolf et al., 2019) under their respective licenses: GPT-2 (Modified MIT), 1078 Llama 3.2 (Meta Llama 3 Community), Qwen 2.5 1079 (Apache 2.0), and Gemma-2B (Gemma Terms of 1080 Use). While all models were pre-trained on English data, LLama and Qwen models provide ad-1082 ditional multilingual capabilities, namely English, 1083 German, French, Italian, Portuguese, Hindi, Span-1084 ish, and Thai for LLama, and over 10 languages, 1085 including Chinese, English, French, Spanish, Por-1086 tuguese, Russian, Arabic, Japanese, Korean, Viet-1087 namese, Thai, and Indonesian for Qwen. The mod-1088 els employ the following tokenisation approaches: 1089 GPT-2, Byte Pair Encoding (BPE) with a 50,257-1090 token vocabulary, optimised primarily for English, 1091 Llama 3.2 uses SentencePiece-based BPE, combining 100K tokens from Tiktoken3 with 28K addi-1093 tional tokens to enhance multilingual performance, 1094 Qwen 2.5 employs Byte-level BPE, utilising a 1095 151,643-token vocabulary designed for multilin-1096 gual processing, Gemma-2B has a SentencePiece 1097 tokeniser leveraging a 256,000-token vocabulary, making it highly effective for English-based tasks. Each model was fine-tuned on its respective down-1100 1101 stream task following a systematic hyperparameter search to identify optimal configurations. Prior 1102 to fine-tuning, prompt engineering was conducted 1103 to determine well-performing prompts tailored to 1104 each task, ensuring alignment with task-specific 1105 requirements and enhancing the models' ability to 1106 generate accurate and contextually relevant outputs. 1107 The hyperparameter search explored key factors, 1108 including weights for stability regularisation, mu-1109 tual information (MI) regularisation, and the over-1110 all CARMA weight (Equation 7), as well as the 1111 specific layers to which these losses were applied. 1112

For training parameters, the following batch sizes were set in the IDM task: 16 for the Gemma-2B and GPT models, 32 for the Qwen-3B and Llama models, and 64 for the Qwen-0.5B model. For SC, the batch sizes were 16 for the GPT models, Gemma-2B and Llama-3B; 32 for Llama-1B and Qwen-3B; and 64 for Qwen-0.5B. For the number of training epochs, in the IDM, the Gemma and GPT models were trained for two epochs, while all other models were trained for three epochs, whereas all models were trained for two epochs, except Gemma-2B and LLama-1B, which were trained for three epochs for the SC task. The stopping layers for IDM and CARMA were configured as follows: GPT2-S at layer 3, GPT2-L at layer 8, Gemma-2B at layer 10, Llama-1B at layer 7,

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123 1124

1125

1126

1127

1128

Llama-3B at layers 8 (stability) and 12 (MI), Qwen-1129 0.5B at layer 5, and Qwen-3B at layer 10. The SC, 1130 the ending layers, 4 for GPT2-S, 12 for GPT2-L, 10 1131 for Gemma-2B, 7, for LLama 1B, 8, for LLama 3B, 1132 5 for Qwen-0.5B and 7 for Qwen-3B. For CARMA 1133 weight, optimal values varied by model size: 0.4 1134 and 0.5 were most effective for larger models. We 1135 hypothesise that CARMA regularisation exhibits 1136 a weaker effect when lower weights are applied, 1137 particularly in larger architectures where stronger 1138 constraints are needed to stabilise compositional 1139 representations. In IDM, GPT2-L and Gemma per-1140 formed best with a weight of 0.3, GPT2-S with 1141 0.2, Llama-1B with 0.4, and Llama-3B with 0.5. 1142 Qwen models used 0.5 and 0.4 for the 0.5B and 1143 3B variants, respectively. For SC Carma weight, 1144 it was 0.4 for Qwen-0.5B and GPT models, 0.5 1145 for LLama-3B and Qwen-3B, and 0.3 for the rest. 1146 For the ending layer, it was 4 for GPT2-S, 12 for 1147 GPT2-L, 10 for Gemma-2B, 7 for LLama-1B, 8 for 1148 LLama-3B, 5 for Qwen-0.5B and 7 for Qwen-3B. 1149

Model	Parameters	Layers	D _{model}	Heads	Activation	MLP Dimension
GPT-2 Small	85M	12	768	12	GELU	3072
GPT-2 Large	708M	36	1280	20	GELU	5120
Gemma-2B	2B	32	4096	16	GELU	8192
LLaMA3.2 1B	1.1B	16	2048	32	SiLU	8192
LLaMA3.2 3B	3.2B	28	3072	24	SiLU	8192
Qwen2.5-0.5B	391M	24	896	14	SiLU	4864
Qwen2.5-3B	3.0B	36	2048	16	SiLU	11008

Table 3: Summary of model architectures. **Parameters**: total number of trainable parameters; **Layers**: total number of transformer layers; D_{model} : size of word embeddings and hidden states; **Heads**: number of self-attention heads; **Activation**: activation function used in feedforward layers; **MLP Dimension**: dimensionality of the feedforward network.

B.4 Evaluation Metrics

1150

This section details the evaluation metrics used1151in the study, including accuracy, synonym consistency, and performance stability.1152

Accuracyis used as a primary measure of model1154performance and is defined as:1155

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN},$$
 (11) 1150

where TP (true positives) and TN (true negatives)1157denote correctly classified instances, while FP1158(false positives) and FN (false negatives) represent misclassified instances.1160

1161SynonymConsistency(ConsistSyn)1162(ConsistSyn)quantifies a model's ability1163to maintain correct predictions after synonym1164replacement. It is computed as:

$$ConsistSyn = \frac{|Correct After Replacement|}{|Correct Before Replacement|} \times 100,$$
(12)

1166where Correct After Replacement refers to the1167number of correct predictions following synonym1168substitution, and Correct Before Replacement de-1169notes the number of correct predictions before sub-1170stitution. The reported results are the averaged1171ConsistSyn across ($N \ge 5$) runs.

1165

1176

1184

1185 1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1172Coefficient of Variation (CV)The coefficient of1173variation (CV) measures the stability of model per-1174formance across multiple runs, with lower values1175indicating greater consistency. It is defined as:

$$CV = \frac{\sigma}{\mu},$$
 (13)

1177where σ represents the standard deviation of model1178performance across runs, and μ denotes the mean1179performance.

1180Normalised Improvement (NI)Normalised Improvement (NI)1181provement (NI) evaluates the relative gain in consistency introduced by a model over a baseline model.1182It is calculated as:

$$NI = \frac{ConsistSyn_{CARMA} - ConsistSyn_{baseline}}{ConsistSyn_{baseline}} \times 100.$$
(14)

This metric captures the percentage improvement in synonym consistency due to a model variant compared to the baseline model.

C Comprehensive Explanation of Evaluation Interventions

C.1 Constituent-Aware Pooling (CAP) Formalisation

Constituent-Aware Pooling (CAP) Formalisation is a method proposed in (Aljaafari et al., 2024) to systematically assess compositional generalisation via aggregating token-level activations into higherlevel semantic representation. Below is a detailed explanation and formalisation of CAP.

1198Overview. CAP aggregates model activations1199at any chosen constituency level (e.g. tokens to1200words), enabling the analysis of compositional de-1201pendencies. The key steps involved are:

- Input Representations: For a given input sequence $X = [x_1, x_2, ..., x_n]$, the model produces inner states $H = [h_1, h_2, ..., h_n]$ at a specific layer. 1203
- Grouping Constituents: Using syntactic parsers such as Benepar (Kitaev et al., 2019; Kitaev and Klein, 2018), or by inversing the model tokeniser function, the sequence is segmented into constituents $C = [c_1, c_2, \ldots, c_m]$, where each c_i represents a phrase or syntactic unit. For the experiments presented in the paper, tokens were grouped into words to form the smallest linguistic units.
- **Pooling Operations:** For each constituent c_i , the corresponding activations $\{h_j | x_j \in c_i\}$ are aggregated into a single representation r_i using a pooling function:

$$r_i = \alpha(\{h_j | x_j \in c_i\})$$
122

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1221

1228

1229

CAP supports three pooling functions:

Maximum pooling: Selects the highest
 activation values as:
 1223

$$\alpha(\{h_j | x_j \in c_i\}) = \max(\{h_j | x_j \in c_i\}),$$
1224

Mean pooling: Computes the average of activation values as:

$$\alpha(\{h_j | x_j \in c_i\}) = \frac{1}{|c_i|} \sum_{j \in c_i} \{h_j | x_j \in c_i\},$$
 1223

- **Sum pooling:** Accumulates activation values as:

$$\alpha(\{h_j | x_j \in c_i\}) = \sum_{j \in c_i} \{h_j | x_j \in c_i\}.$$
 1230

• Updating Representations: The pooled representations $R = [r_1, r_2, \dots, r_m]$ replace the original activations H for further processing.

Evaluation. The impact of CAP is evaluated by 1234 comparing task-specific performance metrics (e.g., 1235 accuracy, F1 score) of models before and after CAP 1236 is applied. This allows for a direct assessment of 1237 how CAP affects compositionality and task per-1238 formance. This paper utilises the word-level CAP, 1239 pooling related token representation to their corre-1240 sponding words. 1241



Figure 7: SC Performance Across Models Under CAP

C.2 Synonym Replacement

A multi-step approach was adopted to ensure re-1243 liable synonym replacements. First, preprocess-1244 ing was applied to filter out words that were un-1245 likely to produce meaningful replacements. Specif-1246 ically, words belonging to NLTK's predefined stopwords list or shorter than two characters were ex-1248 cluded from consideration. The remaining words 1249 were tagged with their part-of-speech (POS) us-1250 ing spaCy's (Honnibal et al., 2020) POS tagger. 1251 Additionally, the sentiment of each word was de-1252 termined using TextBlob (Loria et al.) to ensure 1253 that replacements preserved the semantic tone of 1254 the original text. Next, a synonym vocabulary was 1255 constructed using words extracted from spaCy's 1256 en_core_web_md language model. This vocabulary was filtered to include only alphabetic common 1258 words with high probability scores (greater than -15 in our case), as determined by spaCy's word fre-1260 1261 quency data, while stopwords and rare terms were excluded. This step ensured that the vocabulary consisted of meaningful and contextually appropri-1263 ate words for replacement. For each target word, a list of synonym candidates was generated by it-1265

erating over the constructed vocabulary. The top n candidates were selected based on their semantic similarity to the original word, measured using spaCy's word vectors. Synonyms with high similarity scores and alignment in POS were prioritised to maintain grammatical and contextual coherence in the text.

1266

1267

1268

1269

1270

1271

1272

1273

1275

1276

1277

1278

1279

1280

1281

1282

1283

D InfoNCE for Mutual Information Estimation

Mutual information (MI) quantifies the shared information between two variables X and Y. CARMA leverages MI maximisation to capture dependencies between tokens effectively, thereby enhancing compositional generalisation in LLMs. Specifically, CARMA uses MI, denoted as I(X;Y), to reinforce token-level interactions critical for compositionality. However, direct computation of MI is challenging in practice.

To address this challenge, a variant of InfoNCE1284is employed to estimate MI and approximate these1285dependencies efficiently. Given an anchor token1286hidden state h_i , we construct a corresponding pos-1287itive set H, which contains tokens hidden states1288

Model	Ver.	Task	Int.	CS	CV
GPT2-S	CARMA	IDM	25%	49.17	0.025
	FT	IDM	25%	50.89	0.017
	Org	IDM	25%	52.46	0.044
	CARMA	IDM	40%	35.90	0.0542
	FT	IDM	40%	37.16	0.0628
	Org	IDM	40%	37.20	0.1223
	CARMA	IDM	25%	56.31	0.0164
	FT	IDM	25%	56.95	0.0311
GPT2-L	Org	IDM	25%	51.10	0.1175
011212	CARMA	IDM	40%	43.56	0.0485
	FT	IDM	40%	43.97	0.0459
	Org	IDM	40%	34.68	0.0895
	CARMA	IDM	25%	56.70	0.023
	FT	IDM	25%	57.42	0.030
Gemma-2B	Org	IDM	25%	49.47	0.031
Ocimia-2D	CARMA	IDM	40%	0.4236	0.0174
	FT	IDM	40%	0.4498	0.0249
	Org	IDM	40%	0.3576	0.0480
	CARMA	IDM	25%	58.40	0.0400
	FT	IDM	25%	57.86	0.0385
I lama-1B	Org	IDM	25%	47.55	0.0503
Liama-1D	CARMA	IDM	40%	47.07	0.0476
	FT	IDM	40%	46.75	0.0455
	Org	IDM	40%	33.49	0.0391
	CARMA	IDM	25%	56.98	0.0286
	FT	IDM	25%	54.57	0.0191
Owen-0.5B	Org	IDM	25%	46.84	0.0684
Qweii-0.5D	CARMA	IDM	40%	40.55	0.0397
	FT	IDM	40%	39.69	0.0491
	Org	IDM	40%	32.98	0.0938
	CARMA	IDM	25%	62.00	0.0225
	FT	IDM	25%	61.79	0.0279
Owen_3B	Org	IDM	25%	49.37	0.0441
Qweii-5B	CARMA	IDM	40%	45.05	0.0400
	FT	IDM	40%	45.74	0.0551
	Org	IDM	40%	31.95	0.0688
	CARMA	IDM	25%	62.86	0.015
	FT	IDM	25%	62.22	0.029
Llama-3R	Org	IDM	25%	52.47	0.035
Liama-5D	CARMA	IDM	40%	49.05	0.0297
	FT	IDM	40%	48.31	0.0191
	Org	IDM	40%	36.95	0.0458

Table 4: Model performance (25% and 40% synonym intervention) on the IDM task. Ver.: Version; Int.: Intervention rate; CS: ConsistSyn (%); CV: Coefficient of Variation. Best values in bold.

semantically or syntactically related to h_i . Additionally, we define \mathcal{N} as the set of negative examples consisting of unrelated tokens hidden states.

The InfoNCE objective provides a practical lower bound on I(X;Y) (Oord et al., 2018), as follows:

$$I(X;Y) \ge \mathbb{E}\left[\log \frac{\sum_{h_j \in \mathbf{H}} f(h_i, h_j)}{\sum_{h_j \in \mathbf{H}} f(h_i, h_j) + \sum_{h_k \in \mathcal{N}} f(h_i, h_k)}\right]$$
(15)

where $f(h_i, h_j) = \exp(\sin(h_i, h_j)/\tau)$ is a scaled similarity function, and τ is a temperature parameter. This adaptation of InfoNCE introduces tokenspecific interactions within the layer-wise structure of LLMs, ensuring that dependencies are captured across layers. By maximising mutual information, CARMA aligns the optimisation direction to enhance compositional structures. 1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1312

1313

1314

To extend this approach across layers, the final CARMA MI loss is computed as:

$$\mathcal{L}_{\mathrm{MI}} = -\frac{1}{N} \sum_{i=1}^{N} \left(\log \sum_{\substack{h_j \in \mathbf{H} \\ j \neq i}} \exp\left(\frac{\sin(h_i, h_j)}{\tau}\right) - \log\left(\sum_{\substack{h_j \in \mathbf{H} \\ j \neq i}} \exp\left(\frac{\sin(h_i, h_j)}{\tau}\right) + \sum_{h_k \in \mathcal{N}} \exp\left(\frac{\sin(h_i, h_k)}{\tau}\right) \right) \right),$$
(16)

where h_i is the anchor token, $h_j \in \mathbf{H}$ are positive examples related to $h_i, h_k \in \mathcal{N}$ are negative examples, N is the number of anchors, and $sim(h_i, h_j)$ is a similarity function. The negative sign ensures that MI is maximised during optimisation. Without this negative sign, the objective would incorrectly minimise MI, thereby hindering CG enhancement.

E Extended results

Figures 6 and 7, and Tables 4 and 5 provide addi-1315 tional results for models' performance comparison 1316 under CAP and synonym interventions. CARMA 1317 models show a clear advantage over all models and 1318 tasks. However, the gain is clearer in the IDM case, 1319 where more intricate features and compositional-1320 ity generalisation are required. It is also observed 1321 that the performance of the FT and CARMA mod-1322 els demonstrates similar curves or trends. Given 1323 this observation, we argue that CARMA's improve-1324 ments stem from its learning objectives, which 1325 align closely with cross-entropy loss while explic-1326 itly addressing intermediate representation stability. 1327 The observed improvements are moderate in some 1328 cases, particularly for SC tasks. This behaviour is 1329 expected due to the limited size of the fine-tuning 1330 datasets compared to the original pretraining data 1331 used for these models. Nevertheless, larger models, such as Llama-3B and Gemma-2B, exhibit more substantial improvements with CARMA, demon-1334 strating its scalability with model capacity. 1335

- 1293
- 1295 1296
- 1298

Model	Ver.	Task	Int.	CS	CV
GPT2-S	CARMA	SC	25%	89.03	0.8903
	FT	SC	25%	89.54	0.8954
	CARMA	SC	40%	84.95	0.0095
	FT	SC	40%	85.07	0.0098
	CARMA	SC	25%	88.58	0.0065
	FT	SC	25%	88.04	0.0082
GPT2-L	CARMA	SC	40%	84.61	0.0072
	FT	SC	40%	84.04	0.0073
	CARMA	SC	25%	84.81	0.0069
	FT	SC	25%	81.67	0.0088
Commo 2P	Org	SC	25%	68.14	0.0076
Gemma-2D	CARMA	SC	40%	81.48	0.0102
	FT	SC	40%	74.29	0.0073
	Org	SC	40%	76.06	0.0136
	CARMA	SC	25%	74.03	0.0069
	FT	SC	25%	75.69	0.0044
Llomo 1D	Org	SC	25%	2.65	0.1239
Liailia-1D	CARMA	SC	40%	71.43	0.0065
	FT	SC	40%	74.31	0.0102
	Org	SC	40%	1.73	0.2245
	CARMA	SC	25%	89.66	0.0037
	FT	SC	25%	89.83	0.0085
Owen 0.5B	Org	SC	25%	59.12	0.0691
Qwell-0.5B	CARMA	SC	40%	86.03	0.0084
	FT	SC	40%	86.31	0.0046
	Org	SC	40%	55.27	0.0429
	CARMA	SC	25%	93.65	0.0061
Qwen-3B	FT	SC	25%	93.85	0.0039
	Org	SC	25%	67.63	0.0227
	CARMA	SC	40%	91.26	0.0050
	FT	SC	40%	91.26	0.0050
	Org	SC	40%	64.05	0.0159
Llomo 2D	CARMA	SC	25%	84.83	0.0056
	FT	SC	25%	85.85	0.0065
	Org	SC	25%	35.21	0.0136
Liama-3D	CARMA	SC	40%	82.89	0.0016
	FT	SC	40%	83.55	0.0067
	Org	SC	40%	32.88	0.0188

Table 5: Model performance (25% and 40% synonym intervention) on the SC task. **Ver.**: Version; **Int.**: Intervention rate; **CS**: ConsistSyn (%); **CV**: Coefficient of Variation. **Best values in bold.**