Guiding Explanation-based NLI through Symbolic Inference Types

Anonymous ACL submission

Abstract

001 This work investigates localised, quasisymbolic inference behaviours in distributional representation spaces by focusing on Explanation-based Natural Language Inference (NLI), where two explanations (premises) are provided to derive a single conclusion. We first establish the connection between natural language and symbolic inferences by characterising quasi-symbolic NLI behaviours, named symbolic inference types. Next, we establish the connection between distributional and sym-011 bolic inferences by formalising the Transformer encoder-decoder NLI model as a rule-based neural NLI model - a quasi-symbolic NLI conceptual framework. We perform extensive experiments which reveal that symbolic inference types can enhance model training and inference 017 dynamics, and deliver localised, symbolic inference control. Based on these findings, we conjecture the different inference behaviours are encoded as functionally separated subspaces 021 in the latent parametric space, as the future direction to probe the composition and generalisation of symbolic inference behaviour in distributional representation spaces.

1 Introduction

026

027

Explanatory sentences (Jansen et al., 2018b) can encode hierarchical, taxonomic, and causal relations between concepts (Gardenfors and Zenker, 2015). By understanding and reasoning over these concepts expressed by explanations, humans can make intricate decisions, which is significant in scientific, cognitive, and AI domains. In this work, we focus on the Explanation-based Natural Language Inference (NLI) task where two explanations (premises) are provided to derive a single conclusion. Within this task, a central challenge involves achieving localised and (quasi-)symbolic inference behaviour. E.g., given the two premises: milk is a kind of liquid and *liquids can flow*, one may derive the conclusion milk can flow by localising and substituting the concept *liquids* with milk.



Figure 1: *Quasi-symbolic NLI representation* conceptual framework. Inference types can be encoded as functional subspaces, which are separated or disentangled in parametric space. Thus, by manipulating the inference types, we can deliver localised, symbolic inference control.

A key question then arises: How can we train current Transformer-based NLI models to learn and generalise this quasi-symbolic behaviour in the distributional representation space? Investigating this question allows us to shorten the gap between deep latent semantics and formal linguistic representations (Gildea and Jurafsky, 2000; Banarescu et al., 2013), integrating the flexibility of distributionalneural models with the properties of linguistically grounded representations, facilitating both interpretability (i.e., compositionality (Dankers et al., 2022; Marcus, 2003)) and generative control. 043

044

045

047

051

056

058

060

061

062

063

064

065

Recent studies have demonstrated that the predicate-argument structure and semantic roles from explanatory sentences (Argument Structure Theory - AST representations) (Jackendoff, 1992) can be effectively represented, localised, and disentangled in the latent space of transformer-based models (Zhang et al., 2024a,c). A particular instance of an AST representation is the Abstract Meaning Representation (AMR) (Banarescu et al., 2013), which represents the relations between semantic variables, allowing us to first establish the connection between natural and symbolic language

inferences. Specifically, we leverage the AMR to
systematically characterise quasi-symbolic inference behaviours, named symbolic inference types,
grounded on AMR symbolic graphs. Using the
explanation-based NLI dataset (EntailmentBank,
Dalvi et al. (2021)), we identify ten categories of
symbolic transformations and provide annotations
for 5,134 premise-conclusion pairs in Section 3.

075

081

083

087

089

095

098

100

101

102

103

104

106

107

108

109

110

111

112

113

114

115

116

117

Next, we establish the connection between distributional and symbolic inferences from the perspective of neural representation spaces (see Section 4). An ideal neuro-symbolic NLI model should demonstrate two core representational capabilities: (i) the capacity to encode and to systematically apply inference rules and (ii) the ability to elicit syntacticsemantic features (Valentino, 2022). Motivated by this, we propose quasi-symbolic NLI representation conceptual framework over a Transformer-based encoder-decoder NLI architecture (Figure 1), in which the symbolic inference types are injected to guide the formation of inference behaviours within the latent parametric space. As for the former, explicit supervision on inference types should align the model's reasoning trajectory with target inference behaviours. By varying different inference types, the model should perform rule-based inference behaviour. With respect to the latter, we introduce a feature space (i.e., abstract sentence bottleneck) in the centre of the encoder-decoder architecture. Ideally, this low-dimensional feature space encodes sufficiently abstract, high-level semantic representations during inference.

> We provide extensive experiments to evaluate both capabilities, including the training and inference (Section 5.1), localised inference control (Section 5.2), and feature representation with explanation inference retrieval task (Section 5.3). Experimental results reveal that the symbolic inference type can assist model training, inference, and deliver localised inference control, indicating the possibility of neural NLI models to learn and generalise the inference rules in the distributional space.

In summary, this work provides a complete initial step in investigating the quasi-symbolic inference over distributional semantic spaces, with the following contributions: (1) We first establish the connection between natural and symbolic language inferences from the perspective of linguistics by systematically characterising quasi-symbolic inference behaviours, named symbolic inference types, grounded on the AST/AMR representations. (2) We establish the distributional-symbolic connection from the perspective of neural representation space by proposing the quasi-symbolic NLI representation conceptual framework where the formation of inference behaviours is guided via our symbolic inference types in the latent space. Experimental results showed that the symbolic inference type supervision can improve model training, inference, and localisation. (3) Based on those findings, we conjecture that different inference types are encoded as functional subspaces which are separated or disentangled in the parametric space. We quantitatively evaluate this hypothesis using Neural Tangent Kernel (NTK) theory in Appendix A. 118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

Interpreting and controlling the NLI process from the perspective of the distributional space is a largely promising approach in NLP. To our knowledge, this is the first study to explore the quasisymbolic NLI behaviour, targetting a more universal NLI control and interpretation, rather than a strict symbolic representation or architectural modification. The experimental pipelines are released¹.

2 Related Work

In this section, we review the related work around two topics: *neuro-symbolic representations* and *semantic control over latent spaces*, to highlight the current research limitation and elucidate the motivation underlying our work.

Neuro-symbolic representations. A longstanding goal in NLP is to blend the representational strengths of neural networks with the interpretability of symbolic systems to build more robust NLI models. Current methods usually inject symbolic behaviour through explicit symbolic representations, including graph (Khashabi et al., 2018; Khot et al., 2017; Jansen et al., 2017; Kalouli et al., 2020; Thayaparan et al., 2021), linear programming (Valentino et al., 2022b; Thayaparan et al., 2024), adopting iterative methods, using sparse encoding mechanisms (Valentino et al., 2020; Lin et al., 2020), synthetic quasi-natural language expression (Clark et al., 2020; Yang and Deng, 2021; Yanaka et al., 2021; Fu and Frank, 2024; Weir et al., 2024), symbolic-refined LLMs (Olausson et al., 2023; Quan et al., 2024), etc. Those studies ignore the underlying neuro-symbolic behaviour in neural representation space.

¹https://anonymous.4open.science/r/Inference_ type-5E07/

From an Explainable AI perspective, many studies have shown that neural networks can encode sparse neural-symbolic concepts without explicit symbolic injection across areas like image embedding (Ren et al., 2022; Deng et al., 2021; Li and Zhang, 2023), word embedding (Ethayarajh et al., 2018; Allen et al., 2019; Ri et al., 2023), contextual embedding (Gurnee et al., 2023; Nanda et al., 2023; Li et al., 2024), and LLM interpretation (Park et al., 2024; Templeton et al., 2024). By understanding the symbolic behaviour within neural networks, their decision-making logic can be better interpreted and controlled (Chen et al., 2024).

165

166

167

168

170

171

172

174

175

176

178

179

181

185

187

188

189

190

192

193

194

195

196

197

198

199

201

204

207

210

211

212

213

214

In this work, we draw on quasi-symbolic NLI objectives within distributional neural models, targetting better controllability and interpretability.

Semantic control over latent spaces. Latent variable models, such as VAE (Kingma and Welling, 2013) and Diffusion (Dhariwal and Nichol, 2021), have shown the capability of symbolic representation, control, and interpretation over the distributional space, which are widely deployed in the NLP domain, such as disentangled representation learning (Zhang et al., 2024a) and style-transfer (Liu et al., 2023a; Gu et al., 2023; Zhang et al., 2024b). Guided by semantic annotation, such as labels (Carvalho et al., 2023) and classifiers (Ho and Salimans, 2022), distinct semantic features can be geometrically separated and composed in the latent space, enhancing localisation and interpretability. However, this concept remains under-explored in the NLI domain. Thus, we propose the quasi-symbolic NLI representation conceptual framework and inference types as an initial step to probe the localised, quasi-symbolic NLI behaviour.

In the next section, we start by defining the symbolic inference types for semantically bridging the natural language and symbolic inferences.

3 Defining Symbolic Inference Types

Valentino et al. (2021) has demonstrated that stepwise explanation-based NLI cannot be directly framed as pure logical reasoning. Explanatory chains, while looking plausible at first inspection, commonly have subtler incompleteness and consistency problems from a logical point of view. Meanwhile, explanatory chains corresponding to definable inference patterns and symbolic operations can be localised over the sentence structure. Motivated by this middle ground between logical representations and lexico-semantic inference patterns, 215 we introduce granular inference types based on ex-216 planatory sentences, using AMR to define the sym-217 bolic operations involved in step-wise inference, 218 linking transformations from premises to conclu-219 sions². Table 1 describes the AMR-grounded infer-220 ence types and examples from the EntailmentBank 221 corpus. Next, we define each lexico-semantic inference type and the corresponding symbolic forms.

P1: a scar on the knee is a kind of scar



Figure 2: AMR argument substitution: the inference behaviour is defined as subgraph substitution.

224

225

226

227

228

229

231

232

233

234

235

236

237

The *substitution* category refers to obtaining a conclusion by replacing a predicate/argument term from one premise with a predicate/argument term from the other premise. Possible variations of this category include (1) *argument (ARG) substitution*, (2) *predicate (PRED) substitution*, and (3) *frame (PRED+ARG) substitution*. In this category, one premise is used to connect two terms which are usually connected by *is a kind of, is a source of*, etc. Conceptualising the AMR representation as a graph, this can be symbolically represented as a subgraph substitution operation over the premise graphs, as illustrated in Figure 2. The *PRED sub-*

²Please note that AMR is not used as a representation mechanism in the proposed architecture, but only to precisely ground these symbolic operations within a well-defined semantic representation structure.

Original type	Symbolic type	Prop.	Example entailment relation
	ARG substitution	19%	P1: a scar on the knee is a kind of scar
	(ARG-SUB)		P2: a scar is an acquired characteristic
			C. a scal of the knee is an acquired characteristic
Substitution	PRED substitution	5%	P2: to contain something can mean to store something
Substitution	(PRED-SUB)	570	C: food stores nutrients and energy for living things
			P1: the formation of diamonds requires intense pressure
	Frame substitution	20%	P2: the pressure is intense deep below earth 's crust
	(FRAME-SUB)		C: the formation of diamonds occurs deep below the crust of the earth
	Conditional frame		P1: if something is renewable then that something is not a fossil
Inference from Rule	insertion/substitution	12%	P2: fuel wood is a renewable resource
	(COND-FRAME)		C: wood is not a fossil fuel
	ARG insertion (ARG-INS) Frame conjunction (FRAME-CONJ)		P1: solar energy comes from the sun
		18%	P2: solar energy is a kind of energy
Further Specification			P3: solar energy is a kind of energy that comes from the sun
or Conjunction		6%	P1: photosynthesis stores energy
			P2: respiration releases energy
			C: photosynthesis stores energy and respiration releases energy
Infer Class	ARG/PRED		P1: rock is a hard material
from Properties	(ADC/DDED CEN)	1%	P2: granite is a hard material
	(ARG/FRED-GEN)		C: granite is a kind of rock
Droporty Inharitonaa	(Property Inheritance)	0 10%	P1. blacktop is made of asphalt concrete
Froperty Inneritance	(ARG-SUB-PROP)	0.4%	C: a blackton has a smooth surface
	(//////////////////////////////////////		an optical telescope requires visible light for human to use
Causal Expression	Causality (IFT)	0.8%	clouds / dusts block visible light
Cuusai Expression	Cuusunty (II I)	0.070	if there is clouds or dusts, then the optical telescope cannot be used
			a shelter can be used for living in by raccoons
Example-based Inference	Example (EXAMPLE)	0.9%	some raccoons live in hollow logs
T	······································		an example of a shelter is a raccoon living in a hollow log

Table 1: Examples of symbolic inference types, with their corresponding abbreviations provided in parentheses and used consistently throughout the paper. The EntailmentBank utilised for this task comprises 5,134 instances, with our annotations covering 84% of the (premises, conclusion) cases. These annotations are planned for public release.

stitution category works in a similar manner, but replacing a predicate term. The two predicates are usually linked by the following patterns: " v_1 is a kind of v_2 ", "to v_1 something means to v_2 something", etc. The frame (PRED+ARG) substitution category combines both previous categories by replacing a frame (predicate subgraph) of one of the premises with one from the other premise.

239

241

242

243

245

247

248

249

252

253

254

256

258

259

261

The *further specification or conjunction* category allows for obtaining a conclusion by joining both premises. It includes (4) *ARG insertion* and (5) *frame conjunction*. For *ARG insertion*, the conclusion is obtained by connecting an argument from one of the premises to a frame of the other. As for *frame conjunction/disjunction*, the conclusion is obtained by joining the premises graphs through a conjunction/disjunction node (*and*) or (*or*).

The *inference from rule* category from Dalvi et al. (2021) encompasses a specific instance of insertion or substitution, identified as (6) *conditional frame insertion/substitution*. In this category, a frame is either inserted or replaced as an argument of a premise, following a conditional pathway present in the other premise. This process is illustrated in

Figure 6 (Appendix B).

The inference type *infer class from properties* has been re-categorised as (7) *ARG or PRED generalisation*, where a new *:domain* relation frame is created if both premise graphs differ by a single predicate/argument term. (8) *Property inheritance*, on the other hand, is a special case of *ARG substitution*, where one of the premises describes a *is made of* relationship between the entity in the other premise and its replacement.

Finally, (9) *Causal Expression* and (10) *Examplebased Inference* categories are defined according to the key lexical characteristic of the conclusion, as systematic AMR transformations which could be applied without rephrasing the underlying explanatory sentences could not be determined. More details about the annotation procedure are provided in Appendix B.

Thus far, we have established a connection between natural and symbolic language inferences through the AMR graph. In the next section, we aim to establish the distributional-symbolic NLI connection from the point of neural representation space. 262

287

289

291

295

296

297

301

302

307

310

311

312

313

314

315

316

320

321

323

325

326

4 Quasi-symbolic NLI Conceptual Framework

In this section, we first formalise the concept of Quasi-symbolic NLI³ and then map it to the practical encoder-decoder architectures.

4.1 Quasi-symbolic NLI Formalisation

In this study, we formalise the concept of "quasisymbolic NLI behaviour" as rule-based reasoning over neural representation, where discrete inference behaviours are implemented through differentiable transformations over continuous neural representations. This is achieved by characterising and manipulating quasi-symbolic inference behaviours, denoted by $\pi \in \Pi$, where Π represents the space of all possible inference rules.

The process involves three key stages: (i) Neural Encoding: The premises p_1 and p_2 are encoded into continuous vector representations $(\vec{p_1} \text{ and } \vec{p_2})$ in a neural space. (ii) Rule-Based Reasoning: The encoded representations are transformed using a reasoning function guided by the inference behaviour π . (iii) Neural Decoding: The resulting vector, \vec{c} , is decoded into a natural language conclusion c. Formally, the process can be described as follows:

1.
$$\overrightarrow{p_1}, \overrightarrow{p_2} = f_{encode}(p_1, p_2)$$

2. $\overrightarrow{c} = f_{reason}(\overrightarrow{p_1}, \overrightarrow{p_2}; \pi)$
3. $c = f_{decode}(\overrightarrow{c})$

Here, f_{encode} , f_{reason} , f_{decode} represent the encoding, reasoning, and decoding functions in a neural NLI model. The injection of π should exhibit two advantages:

1. Training Dynamics: During training, explicit supervision on π aligns the model's reasoning trajectory with target inference behaviours, improving conclusion prediction accuracy.

2. Inference Composition: By varying π during inference, the model can separate the semantics of the premises from the inference behaviour. This enables localised, quasi-symbolic NLI control, allowing for flexible and interpretable reasoning.

4.2 Quasi-symbolic NLI Representation

We focus on encoder-decoder architectures (e.g., T5) due to their inherent separation of reasoning

and decoding phases, which naturally accommodates quasi-symbolic reasoning. However, this framework can also be adapted into decoder-only architecture, where the rules are captured, such as through in-context learning (Liu et al., 2024). From a representational perspective, we propose the concepts of latent rule space and feature space to align with the function of the neuro-symbolic NLI model.

327

328

329

331

332

333

334

335

337

338

339

340

341

342

343

344

345

346

347

348

350

351

352

353

354

355

356

357

358

359

360

361

363

364

365

366

367

368

369

370

371

372

374

Latent rule space. The latent rule space refers to the functional parameter space (i.e., models' weights), which captures the structured, rule-based reasoning behaviours $\pi \in \Pi$. We further propose that rule-based reasoning is primarily materialised in the encoder of the encoder-decoder NLI model:

Proposition: The inference behaviour is instantiated at the encoder and can be controlled by the injection of the associated inference type labels.

Due to the page limitation, we provide a formal proof, illustration, and evaluation in Appendix A.

Latent feature space. The latent feature space refers to the input or output embedding space. To evaluate the feature representation capability, we next describe the methodological framework behind the construction of the abstract sentence representation within T5 (named T5 bottleneck).

As for the encoder's final layer output embedding, we compute dimension-wise mean pooling over token embeddings, followed by a multi-layer perceptron to obtain sentence embeddings. As for the decoder's input embedding, we reconstruct token embeddings via linear projection, feeding them into the decoder's cross-attention mechanism. Here, we only describe the optimal setup. We provide a systematic way to choose the best setup in the Appendix C.

5 Empirical Analysis

The experiment addresses three key questions: Section 5.1: (i) Do symbolic inference types enhance model training and inference performance? Section 5.2: (ii) Can these inference types be used for prescriptive inference control? Section 5.3: (iii) Does the incorporation of a sentence bottleneck contribute to improved feature representation? All experimental details are provided in Appendix C.

5.1 Training and Inference Evaluation

Firstly, we evaluate (i) if symbolic inference types enhance model training and inference performance.

³Please note that our objective is not to propose a new neural-symbolic model architecture; rather, we aim to investigate whether a standard NLI model is capable of exhibiting quasi-symbolic NLI behaviour.

We consider three mechanisms to conditionally inject the symbolic inference types into the latent space, which are described below, where p1, p2, and *con* are the premises and conclusion, respectively, and </s> is a special token for sentence separation: **i.** The inference type as the prefix for the premises at the Encoder: *the inference type is* [type] </s> p1 </s> p2**ii.**The inference type asthe prefix for the conclusion in the Decoder: <math></s>*the inference type is* [type]. *con* **iii.** The inference type at the end of the conclusion in the Decoder: </s> *con. the inference type is* [type].

Training dynamics. We first evaluate generative performance after training based on three metrics: BLEURT (Sellam et al., 2020), BLEU (Papineni et al., 2002), and cosine similarity against sentenceT5 (Ni et al., 2021). By comparing the predicted and golden conclusions, we can fairly evaluate the accuracy of the NLI performance. For the baseline, we choose the T5, GPT2 (Radford et al., 2019), Qwen2.5 (Qwen et al., 2025), Llama3.2 (Grattafiori et al., 2024), our T5 bottleneck and Optimus (Li et al., 2020) with 768 latent dimensions as testbed. The performances are measured from the Entailment test set.

390 391

396

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

As illustrated in Table 2, across all baseline models, incorporating inference types into the encoder consistently results in improved performance as measured by BLEU, Cosine, and BLEURT metrics, indicating the inference behaviour is instantiated at the encoder (*Proposition*) (finding1). This finding also suggests that during training, explicit supervision on inference types aligns the model's reasoning trajectory with target inference behaviours, improving conclusion prediction accuracy (finding2). A similar observation is reflected in the test loss curve shown in Figure 10.

Furthermore, previous studies have revealed that the LLM evaluation can be consistent with the results obtained by expert human evaluation (Chiang and Lee, 2023; Liu et al., 2023b; Wang et al., 2023; Huang et al., 2023). Thus, we also conduct a quantitative analysis to measure whether the generated conclusion contradicts the premises through LLM evaluators, including ChatGPT40 as the baseline and GPT40-mini for comparison. Table 3 indicates that EP can consistently result in improved LLM agreement scores. A qualitative evaluation based on the manual check is presented in appendix D (Tables 14 and 15).

Baseline	INJ	BLEU	Cosine	BLEURT
seq2seq1	M: enco	oder-decod	er architect	ture
Т5	DE	0.55	0.96	0.30
original	DP	0.59	0.96	0.34
(amall)	EP	0.65	0.97	0.45
(sman)	NO	0.54	0.96	0.22
Τ5	DE	0.46	0.96	0.23
original	DP	0.53	0.96	0.25
(hara)	EP	0.61	0.97	0.39
(base)	NO	0.57	0.96	0.33
Τ5	DE	0.60	0.97	0.46
original	DP	0.64	0.97	0.44
(lana)	EP	0.67	0.97	0.50
(large)	NO	0.57	0.96	0.31
Causa	lLM: de	coder only	architectur	re
CDT2(-1)	DP	0.28	0.91	-0.90
GP12(XI)	NO	0.27	0.90	-0.97
0 2 5 (0 5 D)	DP	0.65	0.97	0.48
Qwen2.5(0.5B)	NO	0.63	0.97	0.45
L12 2(1D)	DP	0.63	0.97	0.46
Liama5.2(IB)	NO	0.60	0.96	0.42
seq2s	eqLM w	ith sentence	e bottleneck	k
Τ5	DE	0.35	0.91	-0.15
1.J bottleneck	DP	0.39	0.91	-0.13
(base)	EP	0.42	0.92	-0.07
(base)	NO	0.35	0.91	-0.20
	DE	0.26	0.80	-1.11
Optimus	DP	0.25	0.79	-1.14
(BERT-GPT2)	EP	0.09	0.74	-1.17
	NO	0.07	0.74	-1.20

Table 2: Quantitative evaluation on testset, where best results are highlighted in **bold**. Specification for abbreviation. INJ: ways for injecting the information of inference types into the model, it includes DE: decoder end, DP: decoder prefix, EP: encoder prefix, NO: no inference type.

Baseline	INJ	ChatGPT4o	GPT4o-mini
Т5	DE	0.85	0.83
original	DP	0.86	0.83
(lana)	EP	0.91	0.85
(large)	NO	0.84	0.82

Table 3: Agreement scores for the quantitative analysis using LLMs on the test set. We also provide a qualitative manual evaluation in appendix D (Tables 14 and 15), with the prompt being provided in Table 17.

In-context learning. Next, we quantitatively evaluate the symbolic inference types within incontext learning (ICL) in contemporary large language models (LLMs). As illustrated in Table 4, prompting with inference types can improve the performance of ICL in both seq2seq and causal LLMs. Besides, within the context of causal LLMs, an increase in few-shot examples⁴, improves the performance. This finding indicates that our inference types can be generalised across various 425

426

⁴We randomly sample the examples with the same inference type as the current test example from the training set. We perform ten times and calculate the average for each metric.

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

Baseline	INJ	Num	BLEU	Cosine	BLEURT
Seq2seqLLM: encoder-decoder architecture					•
	Yes	10	0.51	0.97	0.39
CoT-T5 (11b)	Yes	5	0.51	0.97	0.39
(Kim et al., 2023)	Yes	0	0.50	0.97	0.36
	NO	0	0.46	0.96	0.31
	Yes	10	0.51	0.97	0.41
El., T5 (1)	Yes	5	0.53	0.97	0.43
Flan-15 (XXI)	Yes	0	0.50	0.96	0.37
	NO	0	0.48	0.96	0.36
Causall	LM: d	lecoder d	only arch	itecture	
	Yes	10	0.52	0.96	0.40
CDT 2.5 tout - 0125	Yes	5	0.48	0.96	0.35
GP1-5.5-turbo-0125	Yes	0	0.46	0.96	0.31
	NO	0	0.42	0.96	0.33
	Yes	10	0.53	0.97	0.50
CDT 4 0612	Yes	5	0.52	0.97	0.47
GP1-4-0015	Yes	0	0.52	0.97	0.50
	NO	0	0.47	0.96	0.40
	Yes	10	0.54	0.97	0.54
11	Yes	5	0.52	0.97	0.52
nama5-706-8192	Yes	0	0.51	0.97	0.47
	NO	0	0.44	0.96	0.40

Table 4: ICL evaluation of test cases, where worst results are highlighted in **bold**. The prompt is "*performing natural language inference [where the inference type is type, description],* $[p1; p2; c]_{\times num}$. p1, p2, what is the *conclusion?*". *num* is the number of examples. The *description* is based on the definition of inference types in Section 3.

5.2 Quasi-symbolic NLI Evaluation

Secondly, we evaluate (ii) if these inference types can be used for prescriptive inference control.

Qualitative evaluation. We qualitatively evaluate the quasi-symbolic NLI behaviour on the generation of conclusions by systematically intervening on the inference type prior to the encoder. As illustrated in Table 5, we can observe that the associated linguistic properties of the conclusion can be controlled consistently with the inference type modifications and this inference control is independent of the semantics of premises, which indicates that the representation mechanisms can improve inference control with regard to symbolic/lexico-semantic properties (finding4). For example, when the type is ARG substitution (ARG-SUB), the model can generate the blacktop is made of a smooth surface by replacing the argument asphalt concrete with smooth surface. The conclusions are changed to asphalt and blacktop have the same surface when the inference type is the conjunction (FRAME-CONJ). Additional examples are provided in Table 16.

P1: blacktop is made of asphalt concrete P2: asphalt has a smooth surface
ARG-SUB: the blacktop is made of smooth surface
ARG-SUB-PROP: blacktop has a smooth surface
ARG/PRED-GEN: a blacktop is a kind of asphalt
ARG-INS: asphalt concrete blacktop has a smooth
surface
FRAME-CON: asphalt and blacktop have the same
surface
IFT: if the asphalt has a smooth surface then the
blacktop will have a smooth surface

Table 5: Controllable generation over original T5 (base) (ARG-SUB: argument substitution, ARG/PRED-GEN: argument/predicate generalisation. ARG-SUB-PROP: property inheritance. ARG-INS: argument insertion, FRAME-CON: frame conjunction, IFT: casual expression.). The example of the T5 bottleneck is provided in Table 12 (Appendix D).

Quantitative analysis. Next, we perform an automated quantitative analysis using LLMs, including ChatGPT4o and GPT4o-mini. For each pair of premises in the EntailmentBank test set, we apply various inference types to generate a diverse set of conclusions using the fine-tuned T5 (base) model. We then assess the resulting (premises, conclusion, inference type) tuples based on two criteria: (i) whether the generated conclusion contradicts the premises, and (ii) whether the (premises, conclusion) pair is consistent with the specified inference type. Utilising the prompt detailed in Table 17 (Appendix D), we report the model agreement score for each criterion. As illustrated in Table 6, the T5 (base) model with controlled symbolic inference types achieves consistency and alignment scores exceeding 60% for both evaluated dimensions.

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

Evaluators	consistency	alignment
ChatGPT4o	0.67	0.63
GPT4o-mini	0.65	0.62

Table 6: Automated quantitative analysis scores.

5.3 Latent Feature Space Evaluation

Finally, we evaluate (iii) whether the incorpora-
tion of feature space (i.e., abstract sentence bot-
tleneck) contributes to improved feature, concept478representation in the NLI task. We especially select
the VAE baselines due to their analogous encoder-
bottleneck-decoder architecture, wherein the com-
pressed sentence bottleneck captures high-level,478

generalised semantics (concepts) (Mercatali and
Freitas, 2021; Zhang et al., 2024a). This structural
similarity is essential for facilitating human-like
inference and cognition (LCM team, 2024).

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

Explanation-based NLI. We first evaluate the NLI performance of different baselines on the Entailment test set. A more effective feature space can enhance generation performance (Zhang et al., 2024c). Consequently, the same generation-related metrics can be applied to evaluate the quality of the feature space.

The baseline includes the state-of-the-art Transformer VAE model: Optimus (Li et al., 2020) and Della (Hu et al., 2022) with two different sentence dimensions (32 and 768). Table 7 illustrates that the T5 bottleneck can outperform all baselines on generation-related metrics, indicating better abstract sentence representations are learned to guide the decoding process.

Baseline	BLEU	Cosine	BLEURT
Optimus(32)	0.07	0.74	-1.20
Optimus(768)	0.08	0.74	-1.21
DELLA(32)	0.08	0.85	-1.23
DELLA(768)	0.09	0.87	-1.09
T5 bottleneck	0.35	0.91	-0.20

Table 7: Comparison of different baselines on EntailmentBank testset, T5 bottleneck has 768 dimensions.

Explanation inference retrieval. We next eval-504 505 uate the abstract sentence embedding using as an associated explanation retrieval task (explanationregeneration - i.e. retrieving the associated explanatory facts relevant to a claim) (Valentino et al., 508 2022a). Given a question-and-answer pair, it reconstructs the entailment tree by searching the ex-510 planations from a fact bank (i.e., WorldTree (Jansen 511 et al., 2018a)) in an iterative fashion using a dense 512 sentence encoder. In this framework, we can re-513 place the sentence embeddings with the proposed 514 T5 bottleneck baseline to evaluate its abstract sen-515 tence embeddings. We compare the T5 bottleneck 516 with abstract sentence representation baselines: Optimus and five LSTM VAEs, and evaluate them via 518 mean average precision (MAP). As illustrated in 519 Table 8, the T5 bottleneck outperforms all base-520 lines, indicating that it can deliver a better abstract 521 representation of explanatory sentences and entailment relations in a retrieval setting (finding5).

depth	t=1	t=2	t=3	t=4
DAE(Vincent et al., 2008)	30.27	31.74	30.65	30.74
AAE(Makhzani et al., 2015)	29.13	30.47	29.33	29.14
LAAE(Rubenstein et al., 2018)	19.13	20.86	18.32	18.01
DAAE(Shen et al., 2020)	13.16	15.42	14.30	13.97
β -VAE(Higgins et al., 2016)	10.03	10.07	10.05	10.05
Optimus(768)	28.21	29.35	28.35	28.27
T5 bottleneck(768)	34.47	35.28	34.50	34.47

Table 8: Explanatory inference retrieval task where t represents the depth of entailment tree.

524

525

526

527

528

529

530

531

532

533

534

535

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

Visualisation. Finally, we visualise the sentence space to evaluate whether inference rules exhibit separability within the latent space. We jointly train the latent space with a linear classifier to predict the inference type categories. As shown in Figure 3, our results indicate that inference types can be partially clustered and separated within this latent space, suggesting the feasibility of rule-based learning through appropriate optimisation strategies (Xie et al., 2025) or architectural designs, such as disentangling rules from lexical semantics.



Figure 3: PCA visualisation: inference types cluster and separation, where left: EP, right: NO.

6 Conclusion

This study serves as a foundational step in exploring the quasi-symbolic NLI behaviour within distributional semantic spaces. We establish the connection between natural and symbolic language inferences by characterising quasi-symbolic inference behaviours based on AMR graphs. Then, we propose the quasi-symbolic NLI representation framework. Experimental results reveal that integrating symbolic inference types enhances training dynamics, inference precision, and explanation retrieval, suggesting the potential for neuro-symbolic NLI. Based on these findings, we hypothesise that distinct inference types can be represented as separated functional subspaces within the parametric space. In future research, we will further examine this hypothesis over different reasoning tasks, such as math reasoning, targetting an explainable and controllable neuro-symbolic NLI model.

Limitations

554

Automatic NLI evaluation. In the domain of 555 LLM automatic evaluation, the prevailing strategy 556 is to select the most advanced LLM as the auto-557 matic evaluator (Chiang and Lee, 2023; Liu et al., 558 2023b; Wang et al., 2023; Huang et al., 2023). We perform a quantitative analysis of the inference 560 consistency in the deductive reasoning process of 562 LLMs, such as ChatGPT-40. However, this assessment may be unreliable due to the inherent limitations of LLMs in logical reasoning. Human evaluation presents a potential alternative, yet the 565 rigorous design of a protocol to systematically verify the logicality of NLI remains an under-explored area in this field. Although we perform a qualitative manual check for LLM evaluation in Table 14 and 569 15, this assessment is not systematic or rigorously 570 structured. A promising direction for improving automatic NLI evaluation is the integration of symbolic theorem provers with LLMs.

Mechanism analysis. This study empirically ex-574 plores quasi-symbolic inference behaviours within 575 distributional semantic spaces. Our findings indicate that symbolic inference types can enhance model training, facilitate inference processes, and enable localised inference control. We hypothesise that quasi-symbolic inference behaviour arises from the geometrical separation of inference types within the parametric space and provide a quantitative evaluation in Appendix A. Future research 583 will further evaluate this hypothesis from different 584 geometric perspectives, such as latent arithmetic, 585 disentanglement, etc., with the target of better com-586 position, arithmetic, generalisation, and interpretation in the neuro-symbolic NLI domain. 588

References

593

594

595

596

598

599

601

602

- Carl Allen, Ivana Balazevic, and Timothy Hospedales. 2019. What the vec? towards probabilistically grounded embeddings. *Advances in neural information processing systems*, 32.
- Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. Abstract meaning representation for sembanking. In *Proceedings of the 7th linguistic annotation workshop and interoperability with discourse*, pages 178–186.
- Danilo S. Carvalho, Yingji Zhang, Giangiacomo Mercatali, and Andre Freitas. 2023. Learning disentangled representations for natural language definitions.

In Findings of the European chapter of Association for Computational Linguistics (Findings of EACL). Association for Computational Linguistics. 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

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

- Lu Chen, Yuxuan Huang, Yixing Li, Yaohui Jin, Shuai Zhao, Zilong Zheng, and Quanshi Zhang. 2024. Alignment between the decision-making logic of llms and human cognition: A case study on legal llms. *arXiv preprint arXiv:2410.09083*.
- Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. *arXiv preprint arXiv:2002.05867*.
- Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan Xie, Hannah Smith, Leighanna Pipatanangkura, and Peter Clark. 2021. Explaining answers with entailment trees. *arXiv preprint arXiv:2104.08661*.
- Marco Damonte, Shay B. Cohen, and Giorgio Satta. 2017. An incremental parser for abstract meaning representation. In *Proceedings of EACL*.
- Verna Dankers, Elia Bruni, and Dieuwke Hupkes. 2022. 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.
- Huiqi Deng, Qihan Ren, Hao Zhang, and Quanshi Zhang. 2021. Discovering and explaining the representation bottleneck of dnns. *arXiv preprint arXiv:2111.06236*.
- Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794.
- Kawin Ethayarajh, David Duvenaud, and Graeme Hirst. 2018. Towards understanding linear word analogies. *arXiv preprint arXiv:1810.04882*.
- Le Fang, Tao Zeng, Chaochun Liu, Liefeng Bo, Wen Dong, and Changyou Chen. 2021. Transformerbased conditional variational autoencoder for controllable story generation. *arXiv preprint arXiv:2101.00828*.
- Xiyan Fu and Anette Frank. 2024. Exploring continual learning of compositional generalization in nli. *arXiv preprint arXiv:2403.04400*.
- Peter Gardenfors and Frank Zenker. 2015. Applications of conceptual spaces: the case for geometric knowledge representation.

Daniel Gildea and Daniel Jurafsky. 2000. Automatic labeling of semantic roles. In *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, ACL '00, page 512–520, USA. Association for Computational Linguistics.

657

661

670

671

672

674

675

683

687

702

704

705

711

714

716

719

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, 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 Allonsius, 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 Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan 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, 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, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Ki720

721

722

723

724

727

728

729

730

731

732

733

734

735

738

740

741

742

743

745

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd of models.

786

794

803

804

810

811

813

814

815

816

817

818

819

820

822

823

825

827

829

830

831

832

833

834

835

837

838

841

843

- Yuxuan Gu, Xiaocheng Feng, Sicheng Ma, Lingyuan Zhang, Heng Gong, Weihong Zhong, and Bing Qin. 2023. Controllable text generation via probability density estimation in the latent space. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 12590–12616, Toronto, Canada. Association for Computational Linguistics.
- Wes Gurnee, Neel Nanda, Matthew Pauly, Kather-

ine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*.

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

897

898

- Irina Higgins, Loïc Matthey, Arka Pal, Christopher P. Burgess, Xavier Glorot, Matthew M. Botvinick, Shakir Mohamed, and Alexander Lerchner. 2016. beta-vae: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*.
- Jonathan Ho and Tim Salimans. 2022. Classifierfree diffusion guidance. *arXiv preprint arXiv:2207.12598*.
- Jinyi Hu, Xiaoyuan Yi, Wenhao Li, Maosong Sun, and Xing Xie. 2022. Fuse it more deeply! a variational transformer with layer-wise latent variable inference for text generation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 697–716, Seattle, United States. Association for Computational Linguistics.
- Fan Huang, Haewoon Kwak, and Jisun An. 2023. Is chatgpt better than human annotators? potential and limitations of chatgpt in explaining implicit hate speech. In *Companion Proceedings of the ACM Web Conference 2023*, WWW '23, page 294–297. ACM.
- Ray S Jackendoff. 1992. *Semantic structures*, volume 18. MIT press.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. 2018. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, volume 31.
- Peter Jansen, Rebecca Sharp, Mihai Surdeanu, and Peter Clark. 2017. Framing qa as building and ranking intersentence answer justifications. *Computational Linguistics*, 43(2):407–449.
- Peter Jansen, Elizabeth Wainwright, Steven Marmorstein, and Clayton Morrison. 2018a. WorldTree: A corpus of explanation graphs for elementary science questions supporting multi-hop inference. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC* 2018), Miyazaki, Japan. European Language Resources Association (ELRA).
- Peter A Jansen, Elizabeth Wainwright, Steven Marmorstein, and Clayton T Morrison. 2018b. Worldtree: A corpus of explanation graphs for elementary science questions supporting multi-hop inference. *arXiv preprint arXiv:1802.03052*.
- Aikaterini-Lida Kalouli, Richard Crouch, and Valeria de Paiva. 2020. Hy-NLI: a hybrid system for natural language inference. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5235–5249, Barcelona, Spain (Online). International Committee on Computational Linguistics.

- 900 901 902 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 926 927 928
- 930 931 932 933
- 934 935 939
- 940 941
- 945
- 946 947
- 948 949

951

952

955

- Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and Dan Roth. 2018. Question answering as global reasoning over semantic abstractions. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 32.
- Tushar Khot, Ashish Sabharwal, and Peter Clark. 2017. Answering complex questions using open information extraction. arXiv preprint arXiv:1704.05572.
- Seungone Kim, Se June Joo, Doyoung Kim, Joel Jang, Seonghyeon Ye, Jamin Shin, and Minjoon Seo. 2023. The cot collection: Improving zero-shot and few-shot learning of language models via chain-of-thought fine-tuning. arXiv preprint arXiv:2305.14045.
- Diederik P Kingma and Max Welling. 2013. Autoencoding variational bayes. arXiv preprint arXiv:1312.6114.
- Paul-Ambroise Duquenne Maha Elbayad Artyom Kozhevnikov Belen Alastruey Pierre Andrews Mariano Coria Guillaume Couairon Marta R. Costajussà David Dale Hady Elsahar Kevin Heffernan João Maria Janeiro Tuan Tran Christophe Ropers Eduardo Sánchez Robin San Roman Alexandre Mourachko Safiyyah Saleem Holger Schwenk LCM team, Loïc Barrault. 2024. Large Concept Models: Language modeling in a sentence representation space.
- Chunyuan Li, Xiang Gao, Yuan Li, Baolin Peng, Xiujun Li, Yizhe Zhang, and Jianfeng Gao. 2020. Optimus: Organizing sentences via pre-trained modeling of a latent space. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4678-4699.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024. Inferencetime intervention: Eliciting truthful answers from a language model. Advances in Neural Information Processing Systems, 36.
- Mingjie Li and Quanshi Zhang. 2023. Does a neural network really encode symbolic concepts? In International Conference on Machine Learning, pages 20452-20469. PMLR.
- Bill Yuchen Lin, Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Xiang Ren, and William W Cohen. 2020. Differentiable open-ended commonsense reasoning. arXiv preprint arXiv:2010.14439.
- Guangyi Liu, Zeyu Feng, Yuan Gao, Zichao Yang, Xiaodan Liang, Junwei Bao, Xiaodong He, Shuguang Cui, Zhen Li, and Zhiting Hu. 2023a. Composable text controls in latent space with ODEs. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 16543–16570, Singapore. Association for Computational Linguistics.
- Sheng Liu, Haotian Ye, Lei Xing, and James Zou. 2024. In-context vectors: Making in context learning more effective and controllable through latent space steering.

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. G-eval: NLG evaluation using gpt-4 with better human alignment. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 2511–2522, Singapore. Association for Computational Linguistics.

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1002

1003

1004

1005

1006

1007 1008

1009

1010

1012

- Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. 2015. Adversarial autoencoders. arXiv preprint arXiv:1511.05644.
- Gary F Marcus. 2003. The algebraic mind: Integrating connectionism and cognitive science. MIT press.
- Giangiacomo Mercatali and André Freitas. 2021. Disentangling generative factors in natural language with discrete variational autoencoders. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 3547–3556, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ivan Montero, Nikolaos Pappas, and Noah A Smith. Sentence bottleneck autoencoders from 2021.transformer language models. arXiv preprint arXiv:2109.00055.
- Neel Nanda, Andrew Lee, and Martin Wattenberg. 2023. Emergent linear representations in world models of self-supervised sequence models. In Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP, pages 16–30, Singapore. Association for Computational Linguistics.
- Jianmo Ni, Gustavo Hernández Ábrego, Noah Constant, Ji Ma, Keith B Hall, Daniel Cer, and Yinfei Yang. 2021. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models. arXiv preprint arXiv:2108.08877.
- Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang, Armando Solar-Lezama, Joshua Tenenbaum, and Roger Levy. 2023. LINC: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 5153–5176, Singapore. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pages 311-318.
- Kiho Park, Yo Joong Choe, and Victor Veitch. 2024. The linear representation hypothesis and the geometry of large language models. In Proceedings of the 41st International Conference on Machine Learning, volume 235 of Proceedings of Machine Learning Research, pages 39643-39666. PMLR.
- Xin Quan, Marco Valentino, Louise A Dennis, and André Freitas. 2024. Verification and refinement of natural language explanations through llm-symbolic theorem proving. arXiv preprint arXiv:2405.01379.

Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 technical report.

1014

1015

1016

1018

1023

1024 1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1063

1064

1066

- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Jie Ren, Mingjie Li, Qirui Chen, Huiqi Deng, and Quanshi Zhang. 2022. Towards axiomatic, hierarchical, and symbolic explanation for deep models.
- Narutatsu Ri, Fei-Tzin Lee, and Nakul Verma. 2023. Contrastive loss is all you need to recover analogies as parallel lines. *arXiv preprint arXiv:2306.08221*.
- Paul K Rubenstein, Bernhard Schoelkopf, and Ilya Tolstikhin. 2018. On the latent space of wasserstein auto-encoders. *arXiv preprint arXiv:1802.03761*.
- Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696*.
- Tianxiao Shen, Jonas Mueller, Regina Barzilay, and Tommi Jaakkola. 2020. Educating text autoencoders: Latent representation guidance via denoising. In *International Conference on Machine Learning*, pages 8719–8729. PMLR.
- Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. 2024. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*.
- Mokanarangan Thayaparan, Marco Valentino, and André Freitas. 2021. Explainable inference over grounding-abstract chains for science questions. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1–12.
- Mokanarangan Thayaparan, Marco Valentino, and André Freitas. 2024. A differentiable integer linear programming solver for explanation-based natural language inference. *arXiv preprint arXiv:2404.02625*.
- Marco Valentino. 2022. Explanation-based scientific natural language inference.

Marco Valentino, Ian Pratt-Hartmann, and André Freitas. 2021. Do natural language explanations represent valid logical arguments? verifying entailment in explainable nli gold standards. 1067

1068

1071

1072

1073

1074

1075

1076

1077

1079

1080

1081

1082

1083

1084

1085

1088

1089

1090

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

- Marco Valentino, Mokanarangan Thayaparan, Deborah Ferreira, and André Freitas. 2022a. Hybrid autoregressive inference for scalable multi-hop explanation regeneration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11403– 11411.
- Marco Valentino, Mokanarangan Thayaparan, and André Freitas. 2020. Explainable natural language reasoning via conceptual unification. *arXiv preprint arXiv:2009.14539*.
- Marco Valentino, Mokanarangan Thayaparan, and André Freitas. 2022b. Case-based abductive natural language inference. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1556–1568.
- Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. 2008. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning*, ICML '08, page 1096–1103, New York, NY, USA. Association for Computing Machinery.
- Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023. Is ChatGPT a good NLG evaluator? a preliminary study. In *Proceedings of the 4th New Frontiers in Summarization Workshop*, pages 1–11, Singapore. Association for Computational Linguistics.
- Nathaniel Weir, Kate Sanders, Orion Weller, Shreya Sharma, Dongwei Jiang, Zhengping Jiang, Bhavana Dalvi Mishra, Oyvind Tafjord, Peter Jansen, Peter Clark, and Benjamin Van Durme. 2024. Enhancing systematic decompositional natural language inference using informal logic. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9458–9482, Miami, Florida, USA. Association for Computational Linguistics.
- Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. 2025. Logic-rl: Unleashing llm reasoning with rule-based reinforcement learning.
- Hitomi Yanaka, Koji Mineshima, and Kentaro Inui. 2021. SyGNS: A systematic generalization testbed based on natural language semantics. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 103–119, Online. Association for Computational Linguistics.
- Kaiyu Yang and Jia Deng. 2021. Learning symbolic rules for reasoning in quasi-natural language. *arXiv* preprint arXiv:2111.12038.

1123Yingji Zhang, Danilo Carvalho, and Andre Freitas.11242024a. Learning disentangled semantic spaces of1125explanations via invertible neural networks. In Pro-1126ceedings of the 62nd Annual Meeting of the Associa-1127tion for Computational Linguistics (Volume 1: Long1128Papers), pages 2113–2134, Bangkok, Thailand. As-1129sociation for Computational Linguistics.

1130

1131

1132

1133

1134

1135

- Yingji Zhang, Danilo Carvalho, Marco Valentino, Ian Pratt-Hartmann, and Andre Freitas. 2024b. Improving semantic control in discrete latent spaces with transformer quantized variational autoencoders. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 1434–1450, St. Julian's, Malta. Association for Computational Linguistics.
- 1137Yingji Zhang, Marco Valentino, Danilo Carvalho, Ian1138Pratt-Hartmann, and Andre Freitas. 2024c. Graph-1139induced syntactic-semantic spaces in transformer-1140based variational AutoEncoders. In Findings of the1141Association for Computational Linguistics: NAACL11422024, pages 474–489, Mexico City, Mexico. Associ-1143ation for Computational Linguistics.

ents for different types:

$$G_{ij}(x,x') := \langle \nabla_{\theta} f_{\theta}(x,\pi_i), \nabla_{\theta} f_{\theta}(x',\pi_j) \rangle \quad (8)$$

If the symbolic inference types π_i and π_j encode 1228 fundamentally different reasoning operations (e.g., 1229 ARG-SUB vs. PRED-SUB), the gradients with 1230 respect to θ for inputs labeled with π_i and those 1231

15

Proposition (Latent Rule Space Hypothesis)

Let \mathcal{M} be a Transformer-based encoder-decoder model parameterised by $\theta = (\theta_{enc}, \theta_{reason}, \theta_{dec})$. Suppose the model performs inference conditioned on a set of symbolic inference types $\Pi =$ $\{\pi_1, \pi_2, \ldots, \pi_n\}$. Then, under supervised training with inference-type annotations, the encoder parameters θ_{enc} induce a parametric structure in which each symbolic inference type π_i corresponds to a distinct subspace $S_{\pi_i} \subseteq \mathbb{R}^{\hat{D}}$ of the encoder representation space.

 $\forall \pi_i, \pi_j \in \Pi, \, \pi_i \neq \pi_j \quad \Rightarrow \quad S_{\pi_i} \cap S_{\pi_i} \approx \emptyset$ (1)

Formal Proof

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164 1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

Step 1: formal setup. Let f_{encode} be the encoder function such that:

$$f_{\text{encode}}: (p_1, p_2, \pi) \mapsto \mathbf{z}_{(p_1, p_2, \pi)} \in \mathcal{Z} \subseteq \mathbb{R}^D \quad (2)$$

where p_1, p_2 are the input premises, $\pi \in \Pi$ is the symbolic inference type, and \mathcal{Z} is the latent representation space of dimension D. The training objective is to minimise the expected loss:

$$\mathcal{L}(\theta_{\text{enc}}) = \mathbb{E}_{(p_1, p_2, \pi, c) \sim \mathcal{D}} \left[-\log p_{\theta}(c \,|\, p_1, p_2, \pi) \right]$$
(3)

where c is the conclusion.

Step 2: NTK interpretation of inference-type subspaces. Each symbolic inference type π is explicitly embedded as part of the model input, for example, as a token prefix (EP). As a result, the model effectively learns a function $f_{\theta}(x, \pi)$, where $x = (p_1, p_2)$ are the premises and π is the symbolic inference type. The function f_{θ} thus jointly depends on both the content of the premises and the nature of the symbolic operation to be performed.

Within the Neural Tangent Kernel (NTK) framework, the similarity between two input examples of the same inference type π is captured by the NTK as follows:

$$\Theta_{\pi}(x, x') = \nabla_{\theta} f_{\theta}(x, \pi)^{\top} \nabla_{\theta} f_{\theta}(x', \pi) \quad (4)$$

where $\nabla_{\theta} f_{\theta}(x, \pi)$ denotes the gradient of the 1182 model output with respect to its parameters, eval-1183 uated at the input (x, π) . This kernel quantifies 1184 how a parameter update from one input-output pair 1185 would affect another pair, conditioned on the shared 1186 inference type. 1187

According to NTK theory (Jacot et al., 2018), in the infinite-width limit, the evolution of the model's predictions under gradient descent training can be described by a linear kernel regression in the RKHS (Reproducing Kernel Hilbert Space) associated with Θ_{π} . Specifically, the prediction at time t, $f_t(x, \pi)$, evolves as:

$$f_t(x,\pi) = f_0(x,\pi) - \Theta_\pi(x,\cdot) \left[\Theta_\pi + \lambda I\right]^{-1} (f_0 - c)$$
(5)

where $f_0(x,\pi)$ is the model's output at initialisation for each training input, λ is a regularisation parameter, and c is the vector of ground truth conclusions.

Crucially, this formulation implies that each symbolic inference type π induces a distinct kernel Θ_{π} , which in turn defines a unique RKHS \mathcal{H}_{π} —that is, a function space within which the model's solutions for inference-type π reside. As the symbolic type π is varied, the structure of the kernel and the corresponding function space changes, reflecting the distinct reasoning behaviours or transformations associated with different inference operations. Thus, the model encodes different symbolic inference patterns in distinct, kernel-induced subspaces.

Step 3: Disjointness via kernel independence. For two different inference types, $\pi_i \neq \pi_j$, we examine the relationship between their corresponding neural tangent kernels (NTKs), Θ_{π_i} and Θ_{π_i} . Specifically, we are interested in the interaction between the parameter gradients induced by inputs associated with different inference types.

Consider two data points x and x', possibly corresponding to different premise pairs. The NTK entry for each type is:

$$\Theta_{\pi_i}(x, x') = \nabla_\theta f_\theta(x, \pi_i)^\top \nabla_\theta f_\theta(x', \pi_i) \quad (6)$$

and

$$\Theta_{\pi_j}(x, x') = \nabla_\theta f_\theta(x, \pi_j)^\top \nabla_\theta f_\theta(x', \pi_j) \quad (7)$$

When considering cross-type similarities, we are interested in the inner product between the gradi-

1194

1188

1189

1190

1191

1192

1193

1195

1196 1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225



Figure 4: Quantitative measuring the separability between different inference-type subspaces in T5 (small) where left: NO, right: EP. We can observe that when injecting inference-type categories into the model during training, the diagonal values exhibit higher values, indicating the inference-type subspaces can be better separated in the parameter space.

labeled with π_j will tend to point in different directions in parameter space. This is because each type imposes a distinct task or transformation pattern on the model, causing it to utilise different portions of its capacity.

1232

1233

1234

1235

1236

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

Under idealised training, where the data for each type is sufficiently distinct and the network has enough capacity, the gradients for one type will have minimal overlap with those of the other. This can be formalised by observing that:

$$\langle \nabla_{\theta} f_{\theta}(x, \pi_i), \nabla_{\theta} f_{\theta}(x', \pi_j) \rangle \approx 0 \qquad \text{for } \pi_i \neq \pi_j$$
(9)

This property implies that the parameter updates driven by examples from different inference types are approximately orthogonal, meaning that training on one type will not interfere with or alter the function learned for the other type. In the language of NTK and kernel regression, this corresponds to the induced RKHS for each type, $\mathcal{H}\pi_i$ and $\mathcal{H}\pi_j$, being approximately disjoint:

$$\mathcal{H}\pi_i \cap \mathcal{H}\pi_j \approx \emptyset \tag{10}$$

Quantitative evaluation. Therefore, by measuring the cosine similarity between gradient vectors associated with different inference types, we can quantify the separability between different inference-type subspaces in the T5 (small) model, comparing settings without (left: NO) and with (right: EP) the injection of inference-type categories during training. As shown in Figure 4, the diagonal values are notably higher in the EP condition, suggesting that incorporating inference-type information during training enhances the separation of inference-type subspaces in the model's parameter space. Thus, NTK theory supports the empirical observation that symbolic inference types induce separated, controllable latent subspaces.

Corollary (Controllability). Therefore, given that each inference type π_i induces a distinct RKHS \mathcal{H}_{π_i} , the inference function becomes selectively controllable:

$$\pi \mapsto f_{\text{reason}}(\mathbf{z}_{(p_1, p_2, \pi)}) \tag{11}$$

1262

1263

1264

1265

1266

1267

1268

1269

1272

1273

1274

1275

1276

1277

1278

1279

1280

1282

1283

1284

1285

1287

1288

1289

1290

1291

Hence, manipulating π explicitly prescribes the symbolic reasoning pattern applied during inference, as we evaluated in Section 5.2. More controlled examples are provided in Table 16.

Feature space separation. Each \mathcal{H}_{π} (the RKHS induced by the NTK for type π) maps to a region of output space. If $\mathcal{H}_{\pi_i} \cap \mathcal{H}_{\pi_j} \approx \emptyset$, the downstream classifier will operate on disjoint regions of the feature space, as visualised in Figure 3.

B Annotation Details

Annotation procedure. Annotation was performed manually for 5134 entailment triples (two premises, one conclusion) from the Entailment-Bank (Dalvi et al., 2021), according to Algorithm 1. Graph subset relations and root matching are relaxed for non-argument (:ARG*, op*) edges, meaning relations such as *:manner* or *:time* can be ignored for this purpose. Two independent annotators with post-graduate level backgrounds in Computational Linguistics were used in this process, on a consensus-based annotation scheme where a first annotator defined the transformations and a second annotator verified and refined the annotation scheme, in two iterations. The annotation of the AMR graph is based on an off-the-shelf parser (Damonte et al., 2017). The descriptions for each inference type category are as follows:

ARG-SUB (Figure 2): the conclusion is obtained by replacing one argument with another argument.

PRED-SUB: the conclusion is obtained by replacing one verb with another verb.

FRAME-SUB: the conclusion is obtained by replacing a frame of one of the premises with one from the other premise.

COND-FRAM (Figure 6): the conclusion is obtained according to the conditional premise with keyword "if".

ARG-INS (Figure 5): the conclusion is obtained by connecting an argument from one of the premises to a frame of the other.

FRAME-CONJ: the conclusion is obtained by using connectives to connect two premises.

ARG/PRED-GEN (Figure 7): a new *:domain* relation frame is created in the conclusion if both premise graphs differ by a single predicate/argument term.

ARG-SUB-PROP (Figure 8): one of the premises describes a "*is made of*" relationship between the entity in the other premise and its replacement.

IFT: the conclusion should be a conditional sentence.

EXAMPLE: the conclusion should contain the keyword "example".

P1: energy comes from food P2: healing requires energy



Figure 5: AMR argument insertion (ARG-INS).

P1: inventing paper allows paper to be used



C: inventing paper might increase the use of paper



Figure 6: AMR conditional frame insertion (COND-FRAME).



Figure 7: AMR argument generalisation (ARG-GEN).

Unknown (UNK) category. In this work, our annotation occupies 84% based on the Entailment-Bank corpus. As for other unknown categories, we do not further specify them, as they either require knowledge outside of the scope of the premises or do not have a consistent symbolic transformation expression. An additional subtype called *premise copy* was included for the cases where the conclusion has the same graph as one of the premises.

1334

1335

1327

1292

1293

1295

1297

1298

1299

1300

1301

1302

1303

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324



Figure 8: AMR argument substitution (property inheritance) (ARG-SUB-PROP).

С **Experimental Details**

C.1 Dataset

1336

1337

1339

1340

1342

1343

1344

1345

1346

1347

Table 9 describes the statistical information of the corpus used in the experiment. For experiments: Section 5.1, 5.2, and 5.3, the Entailment-Bank dataset is split into train 60%, valid 20%, and test 20% sets. For the explanation inference retrieval task in Section 5.3, we follow the same experimental setup provided online.⁵

Corpus	Num data.	Avg. length
WorldTree (Jansen et al., 2018a)	11430	8.65
EntailmentBank (Dalvi et al., 2021)	5134	10.35

Statistics from explanations datasets. Table 9: WorldTree is used in the Explanation Inference Retrieval task.

C.2 T5 Bottleneck Architecture

Figure 9 shows the architecture of the T5 bottleneck for learning latent sentence space. It includes

two stages: sentence embedding and decoder connection. The sentence embedding aims to transform token embeddings into a sentence (single) embedding. Decoder connection aims to connect the encoder and decoder.

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1377

1378

1379

1380

1381

1382

1383

1384

1387

1388

1389

1391

1395

1396

1397

Latent sentence space: While designing the sentence bottleneck, we compare the four most frequently used mechanisms to transform token embeddings into sentence embeddings:

(1) Mean pooling: calculating the mean of each dimension on all token embeddings and feeding the resulting vector into a multi-layer perceptron to obtain the sentence embedding. (2) multi-layer perceptron (MLP): applying an MLP to reduce the dimensionality of token embeddings, and the resulting embeddings are concatenated to form a single sentence embedding: $z = \operatorname{concat} |\operatorname{MLP}_1(x_1); ...; \operatorname{MLP}_T(x_T)|$ where $MLP_i(x_i)$ represents the *i*-th neural network for input representation of token x_i , z is the latent sentence representation, and T is the maximum token length for a sentence. (3) multi-head attention: feeding each token embedding into the multi-head attention and considering the first output embedding as the sentence embedding (Montero et al., 2021): z =MultiHead (XW^q, XW^k, XW^v) [0] where $X = [x_1, ..., x_T]$ and W^q , W^k , and W^v are the weights for learning q, k, v embeddings in selfattention, respectively. (4) Sentence T5: re-loading the pre-trained sentence T5 (S-T5, Ni et al. (2021)).

Conditional generation: Next, we consider four strategies to inject sentence embeddings into the decoder. (1) Cross-attention input embedding (CA Input): reconstructing the token embeddings from a sentence representation and directly feeding them into the cross-attention layers of the decoder: Y =MultiHead $(YW^q, MLP(z)W^k, MLP(z)W^v)$

where \hat{Y} is the reconstruction of decoder input 1385 sequence $Y = [y_1, ..., y_K]$. (2) Cross-attention 1386 KV embedding (CA KV): instead of reconstructing the token embeddings, it consists of directly learning the Key and Value (Hu et al., 2022; Li et al., 2020), which is formalised as $\hat{Y} =$ MultiHead $(YW^q, MLP_k(z), MLP_v(z)),$ where MLP_k and MLP_v are neural layers for learning k v embeddings. (3) Non-cross-attention 1393 input connection (NCA Input): reconstructing 1394 the token embeddings and element-wisely adding them with the input embeddings of the decoder (Fang et al., 2021). (4) Non-cross-attention

⁵https://github.com/ai-systems/hybrid_ autoregressive_inference

output connection (NCA Output): adding the reconstructed token embeddings to the output embedding of the decoder.

Train: architecture					
Decoder Co	nnection	CA Input	CA KV	NCA Input	NCA Output
Sentence Embedding	Pooling	1.41	1.44	1.86	2.42
	MLP	1.71	1.94	2.09	2.62
	MHA	1.51	2.24	2.31	3.03
	S-T5	1.24	1.42	1.81	2.22

Table 10: Comparison of different setups on test loss via cross-entropy (CA: cross-attention, NCA: non-cross-attention), bottom: comparison of different baselines on EntailmentBank testset.



Figure 9: The architectural configuration of T5 bottleneck, it consists of two stages: sentence embedding and decoder connection.

C.3 Implementation Details

Hyper-parameters. 1. Size of Sentence Representation: in this work, we consider 768 as the size of the sentence embedding. Usually, the performance of the model improves as the size increases.
2. Multi-head Attention (MHA): in the experiment, MHA consists of 8 layers, each layer containing 12 heads. The dimensions of Query, Key, and Value are 64 in each head. The dimension of token embedding is 768. Training hyperparameters are: 3. For all models, the max epoch: 40, learning rate: 5e-5. During fine-tuning the T5 bottleneck, we first freeze the pre-trained parameters in the first epoch and fine-tune all parameters for the remaining epochs. 4. All models are trained on a single A6000 GPU device.

Baselines. In the experiment, we implement five1418LSTM-based autoencoders, including denoising1419AE (Vincent et al. (2008), DAE), β -VAE (Hig-1420gins et al., 2016), adversarial AE (Makhzani et al.1421(2015), AAE), label adversarial AE (Rubenstein1422et al. (2018), LAAE), and denoising adversarial



Figure 10: The test loss curve indicates that EP facilitates better convergence, indicating the supervision on inference types aligns the model's reasoning trajectory with target inference behaviours, improving conclusion prediction accuracy.

autoencoder (Shen et al. (2020), DAAE). Their implementation relies on the open-source codebase available at the URL ⁶. As for transformer-based VAEs, we implement Optimus (Li et al., 2020)⁷ and Della (Hu et al., 2022)⁸. All baseline models undergo training and evaluation with the hyper-parameters provided by their respective sources. A latent dimension of 768 is specified to ensure a uniform and equitable comparative analysis.

Metrics. To evaluate the generated conclusions against the reference conclusions, we employ BLEU scores for 1- to 3-gram overlaps and report the average score. Additionally, to assess semantic similarity, we calculate the cosine similarity between the generated and reference conclusions by encoding both using the pretrained Sentence-T5 model⁹ and computing the cosine similarity of their resulting embeddings.

D Complementary Results

Ablation studies. We remove the inference types from the dataset and evaluate the T5 model performance using the same metrics. In this case, we can compare the model performance trained with or without that inference type. From Table 11, we can observe that the baselines (T5 small and base)

text-autoencoders

⁷https://github.com/ChunyuanLI/Optimus

⁸https://github.com/OpenVLG/DELLA

⁹https://huggingface.co/sentence-transformers/ sentence-t5-base

⁶https://github.com/shentianxiao/

achieve higher BLEU and BLEURT scores without the data with ARG-INS, COND-FRAME, and UNK inference type, respectively. This result indicates that the T5 cannot generalize well over those inference types. Also, removing the UNK inference type from data can achieve lower loss and PPL, which indicates that it has a negative impact on model training.

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

Remove	T5	BLEU	BLEURT	Cosine	Loss \downarrow	$\text{PPL}\downarrow$
FRAME-	small	0.50	0.19	0.95	0.95	2.58
SUB	base	0.60	0.33	0.96	0.72	1.95
ADC INS	small	0.54	0.27	0.95	0.82	2.22
AKO-IINS	base	0.63	0.46	0.97	0.64	1.73
FRAME-	small	0.53	0.26	0.96	0.84	2.28
CONJ	base	0.60	0.35	0.96	0.65	1.76
COND-	small	0.55	0.25	0.96	0.88	2.39
FRAME	base	0.59	0.36	0.96	0.69	1.87
UNK	small	0.55	0.23	0.95	<u>0.53</u>	1.44
UNK	base	0.62	0.40	0.96	<u>0.58</u>	1.57
No	small	0.54	0.22	0.96	0.69	2.22
No	base	0.57	0.33	0.96	0.61	1.65

Table 11: Ablation study over inference type (No: no inference types are removed).

More controllable inference examples. We provide more controlled examples based on both the Original T5 and T5 bottleneck in Table 12, 13, and 16. All examples reveal that the inference type can provide quasi-symbolic inference control to language models.

Quasi-symbolic NLI control
P1: a pumpkin contains seeds P2: fruit contains seeds
Original T5: ARG-INS: a fruit in a pumpkin contains seeds FRAME-CONJ: a pumpkin and fruit both contains seeds FRAME-SUB: fruit is a kind of pumpkin
T5 bottleneck: ARG-INS: fruit is a part of pumpkin that contains seeds FRAME-CONJ: a fruit contains seeds FRAME-SUB: a pumpkin is a kind of plant

Table 12: Controlled generation. original T5(base) (top) and T5 bottleneck (bottom).

Qualitative evaluation for LLM evaluators.We conduct a qualitative evaluation through manualinspection. However, this assessment is not systematic or rigorously structured as we discussed in the

Quasi-symbolic NLI control	
P1: eating something has a negative impact on that something P2: some animals eat cacti ARG-INS: some animals have a negative impact on cacti by eating cacti PRED-SUB: some animals may have a negative impact on cacti FRAME-SUB: eating cacti has a negative impact on that cacti	
ARG-INS: some animals have a negative impact on cacti by eating cacti PRED-SUB: animals have a negative impact on cacti FRAME-SUB: eating cacti has a negative impact on that cacti	

Table 13: Controlled generation. original T5(base) (top) and T5 bottleneck (bottom).

Limitations section. Tables 14 and 15 present ex-
amples with discrepancies in scores between Chat-
GPT40 and GPT40-mini, as well as a compari-
son of predictions between encoder prefix injection
(EP) and the absence of inference-type injection
(NO), respectively.1466
1467
1468
1469

1472

1473

1474

1475

From both tables, we observe that ChatGPT40 tends to be more accurate than GPT40-mini and that EP outperforms NO in generating correct predictions.

Premises	Prediction(NO)	Golden	ChatGPT4o	GPT4o- mini	Human Check
p1: the metal on the roof of a car is in contact with air p2: air contains oxygen and wa-	the car roof contains water vapor and oxygen	the metal on the roof of a car is in contact with oxygen and water vapor	0	1	0 (invalid predi- cate "contains")
p1: friction occurs when the student is rubbing his hands to- gether p2: friction causes the tempera-	rubbing your hands together causes the temperature of the ob- ject to increase	friction causes the temperature of student's hands to increase	1	0	1 (replacing "friction" with "rubbing hands together")
p1: a caterpillar is a kind of in- sect	metamorphosis is when a cater- pillar changes from an immature form to an adult form	an example of metamorphosis is when a caterpillar changes from an immature form to an adult form	1	0	l (replacing "insect" with "metamorpho- sis")
insect changes from an imma- ture form to an adult form p1: an increase in water has a positive impact on alligators p2: a flood is caused by an in-	a flood has a positive impact on alligators	a flood has a positive impact on alligators	1	0	1 (exact match)
p1: predators eat prey	predators catching prey requires catching prey	predators must catch prey to eat prey	0	1	0 (fail to do substitution be- tween "eating" and "catching")
p1: a leaf uses chlorophyll to produce carbohydrates p2: carbohydrates are made of	a leaf uses chlorophyll to pro- duce sugars	a leaf uses chlorophyll to pro- duce sugar	0	1	1 (valid infer- ence)
p1: salt is a kind of pure sub- stance	salt and pepper are kinds of sub- stances	salt and pepper are two sub- stances	1	0	1 (valid con- junction both "salt" and "pep- per")
 p2: pepper is a kind of substance p1: some bacteria are good for humans by helping digestion p2: digestion is when stomach 	some bacteria help digestion by breaking down food	some bacteria are good for hu- mans by helping to break down food	0	1	0 (without "good for hu- man")
acid breaks down food p1: a doorbell is a kind of elec- tric device	closing a doorbell causes the doorbell to function	an electric circuit causes a door- bell to function	0	1	0 (incorrect replacing "elec- tric circuit" with "door- bell")
p2: closing an electric circuit causes an electrical device to function					
p1: burning something is a kind of action	burning something is required for something to occur	if something occurs by burn- ing something else, then burning that something else is required for that something to occur	1	0	1 (valid in- ference by replacing "ac- tion" with "burning some- thing")
p2: if something occurs by per- forming an action, then that ac- tion is required for that some- thing to occur					
p1: echolocation can be used for hunting prey	bats must emit sound in order to hunt prey	bats must emit sound in order to hunt their prey	1	0	l (valid in- ference by replacing "echolocate" with "hunting prey")
der to echolocate p1: different solids will have the same physical properties p2: an mixture is formed by two or more substances combined to- gether physically	one solid will form a mixture	different solids that are com- bined will become a mixture	1	0	0 (incorrect "one solid")

Table 14: Qualitative evaluation for examples with discrepancies in scores between ChatGPT40 and GPT40-mini (NO: no inference type injection, 0: invalid, 1: valid). We can observe that the ChatGPT40 tends to be more accurate than GPT40-mini by human check.

Premises	Prediction(NO)	Prediction(EP)	Golden	ChatGPT40	Human Check
p1: the metal on the roof of	the car roof contains water	the car roof is in contact	the metal on the roof of a	NO:0, EP:1	NO:0, EP:1
a car is in contact with air	vapor and oxygen	with oxygen and water va- por	car is in contact with oxygen and water vapor		
p2: air contains oxygen and water vapor					
p1: a beak is used for catch-	ads are used for eating by birds to catch food	a beak is used for eating by	a beak is used for eating food by some birds	NO:0, EP:1	NO:0, EP:1
p2: eating food requires		Some on us	lood by some onds		
p1: predators must catch	animals must catch and eat	animals must catch prey to	some animals must catch	NO:0, EP:1	NO:0, EP:1
prey to eat prey	prey	eat prey	prey to eat		
p2: a predator is a kind of animal					
p1: an adaptation is a kind	an adaptation is something	adaptation is when a living	an adaptation is a kind of	NO:0, EP:1	NO:0, EP:0
of change	a living thing responds to a change in an environment	thing responds to a change in an environment	change in response to a change in an environment		
p2: adapting is when a					
change in an environment					
p1: a doorbell is a kind of	closing a doorbell causes the	closing an electric circuit	an electric circuit causes a	NO:0, EP:1	NO:0, EP:1
electric device	doorbell to function	causes a doorbell to function	doorbell to function		
cuit causes an electrical de-					
vice to function				NO 4 ED 0	
p1: green plants are made of plant cells	a producer is made of plant cells	producers are made of plant cells	producers are made of plant cells	NO:1, EP:0	NO:1, EP:1
p2: green plants are a kind of producer	cons	cons	cons		
p1: the iron nail has rusted	iron nails rusting is when	the iron nail has rusted	a chemical reaction has hap-	NO:0, EP:1	NO:0, EP:0
	the iron nails chemically re- act with water and oxygen to form iron nail		pened on the iron nail		
p2: rusting is when iron					
chemically reacts with wa-					
p1: wood burns	wood burns when intro-	wood chips burn	wood chips burn	NO:0, EP:1	NO:0, EP:1
I	duced to wood	I	I · · ·	,	,
p2: wood chips are made of wood					
p1: some bacteria are good	some bacteria help digestion	some bacteria are good for	some bacteria are good for	NO:0, EP:1	NO:0, EP:1
gestion	by breaking down tood	tion by breaking down food	down food		
p2: digestion is when stom-					
ach acid breaks down food	a lung is an organ directly	a lung is a kind of organ for	lunge can be used for di	NO:0 EP:0	NO-1 EP-1
gan for directly breathing in	for directly breathing in	directly breathing in smoke	rectly breathing in smoke	NO.0, EF.0	NO.1, EF.1
gasses	smoke		, ,		
p2: smoke is a kind of gas	the hot and cold water	the hot agg will decrease in	the water will increase in	NO:0 EP:1	NO:0 EP:1
in temperature	will increase in temperature	temperature and the cold wa- ter will increase in tempera-	temperature and the egg will decrease in temperature	NO.0, EF.1	NO.0, EF.1
p2: the cold water will in-		uit			
crease in temperature					
p1: plant reproduction re- quires pollinating animals	plants reproduction requires bees that carry pollen	a bee can help plant repro- duction by carrying pollen	a bee can help on pollina- tion in plant reproduction by	NO:0, EP:1	NO:0, EP:1
for pollination			carry pollen		
p2. a bee can help on polli- nation by carrying pollen					
p1: a leaf uses chlorophyll	a leaf uses chlorophyll to	a leaf uses chlorophyll to	a leaf uses chlorophyll to	NO:1, EP:1	NO:1, EP:1
to produce carbohydrates	produce sugars	produce sugar	produce sugar		
p2. carbonyurates are made of sugars					

Table 15: Qualitative evaluation for prediction through EP and NO (NO: no inference type, EP: encoder prefix), we can observe that EP outperforms NO in generating correct predictions.

Premises	Inference Type	T5 original
P1: a pumpkin contains seeds	ARG-INS	a fruit in a pumpkin contains seeds
P2: fruit contains seeds	FRAME-CONJ	a pumpkin and fruit both contain seeds
	IFT	if a pumpkin contains fruit then the fruit may contain seeds
	EXAMPLE	fruit is an example of pumpkins being sown
	ARG/PRED-GEN	a pumpkin is a kind of fruit
	AKG-SUB UNK	a numpkin can contain seeds
	FRAME-SUB	fruit is a kind of pumpkin
P1: sunlight is a kind of solar en-	ARG-INS	solar energy is a kind of resource for plants that uses water
ergy		
P2: water and sunlight are re- sources for plants	FRAME-CONJ	water and sunlight are resources for plants and are kinds of solar energy
	UNK	the resources for plants include water and solar energy
	ARG-SUB	water and solar energy are resources for plants
P1: to move something can mean to transfer something	ARG-SUB	flowing can mean to transfer energy
P2: flowing is a kind of movement for energy	INF	if something flows, then that energy will flow
	FRAME-CONJ	moving can transfer energy and mean flowing
	ARG-INS	flowing can be a kind of transfer of energy to another entity
	ARG/PRED-GEN	transferring energy with flowing can be seen as transferring energy
P1: if a pot is exposed to a stove then that pot may become hot	COND-FRAME	the pot will become hot
P2: the pot is exposed to a stove	EXAMPLE	an example of hot pot is a pot exposed to a stove
P1: eating something has a nega- tive impact on that something	FRAME-SUB	eating cacti has a negative impact on that cacti
P2: some animals eat cacti	PRED-SUB	some animals may have a negative impact on cacti
	ARG-INS	some animals have a negative impact on cacti by eating cacti
	EXAMPLE	cooking cacti is an example of a negative impact on a cactus
	INF	if a cactus has a negative impact on an animal, that cactus could be devoured
P1: seeing requires light	ARG-SUB	reading requires light
P2: reading requires seeing	ARG-INS	light is a kind of requirement for reading
	INF EVAMDLE	if light is moving then reading may be taken
	LINK	a light build will be used for reading
	entr	
Premises	Inference Type	T5 bottleneck
P1: a pumpkin contains seeds	ARG-INS	fruit is part of a pumpkin that contains seeds
P2: fruit contains seeds	FRAME-CONJ	a fruit contains seeds
D1, cuplicate is a kind of color on	FRAME-SUB	a pumpkin is a kind of plant
ergy	AKG-IIN5	water is a kind of resource that is used by plants for growin
P2: water and sunlight are re- sources for plants	FRAME-CONJ	plants and water are resources that require water and energy
	ARG-SUB	plants use water and sunlight to produce energy
P1: to move something can mean to transfer something	ARG-SUB	flowing can mean to transfer energy
P2: flowing is a kind of movement for energy	INF	if something flows, then that energy will flow
	FRAME-CONJ	moving can transfer energy and mean flowing
	ARG-INS	flowing can be a kind of transfer of something
	ARG/PRED-GEN	transferring energy with flowing can be seen as transferring energy
P1: if a pot is exposed to a stove then that pot may become hot	COND-FRAME	the pot may become hot
P2: the pot is exposed to a stove	ARG/PRED-GEN	the pot may be a source of heat
P1: eating something has a nega-	FRAME-SUB	eating cacti has a negative impact on that cacti
tive impact on that something		
P2: some animals eat cacti	PRED-SUB	animals have a negative impact on cacti
P2: some animals eat cacti	PRED-SUB ARG-INS	animals have a negative impact on cacti some animals have a negative impact on cacti by eating cacti
P2: some animals eat cacti P1: seeing requires light	PRED-SUB ARG-INS ARG-SUB	animals have a negative impact on cacti some animals have a negative impact on cacti by eating cacti reading requires light
P2: some animals eat cacti P1: seeing requires light P2: reading requires seeing	PRED-SUB ARG-INS ARG-SUB FRAME-CONJ	animals have a negative impact on cacti some animals have a negative impact on cacti by eating cacti reading requires light reading and feeling can both be used

Table 16: controllable NLI via inference type (Top: original T5, bottom: T5 bottleneck).

Algorithm 1 Annotation procedure

1: Find premise P_x most similar to the conclusion C, $P_{\bar{x}}$ being the other premise. 2: $G_{x,\bar{x},C} \leftarrow \text{AMR graph of } P_x, P_{\bar{x}}, C$, respectively. 3: #----- common ARG-SUB, PRED-SUB ------4: if $G_x = G_c$ or $G_{\bar{x}} = G_c$ then type = PREM-COPY # Comment: no reasoning happen. 5: 6: else if P_x and C differ by one word w then # Comment: common ARG(PRED)-SUB. if w is a verb then 7: type = PRED-SUB8: 9: else 10: type = ARG-SUBend if 11: 12: else 13: # - - ------ COND-FRAME, FRAME-SUB, ARG-SUB-PROP -----Get AMR graphs G_1, G_2, G_c for P_1, P_2 and C respectively. $P_x \to G_x$. $14 \cdot$ if \exists :ARG* $(x, a) \in C$ and $a \in P_{\bar{x}}$ then 15: 16: if \exists :condition($root(G_x)$, $root(G_{\bar{x}})$) then # Comment: see Figure 6, two root nodes are connected by :condition edge 17: type = COND-FRAME 18: 19: else if root(a) is a noun then if $root(G_{\bar{x}}) =$ "make-01" and \exists :ARG*($root(G_{\bar{x}})$, a) then 20: # Comment: "make" as a trigger to classify ARG-SUB and property inheritance. 21: type = ARG-SUB-PROP22: else 23: type = ARG-SUB # ARG-SUB that was not caught by the simpler rule on line 10, 24: due to Px differing from C by more than a single word end if 25: 26: else type = FRAME-SUB27: end if 28: ----- Further-specification and Conjunction ------29: # - - - else if $G_x \subset G_c$ and $G_{\bar{x}} \subset G_C$ then 30: 31: type = FRAME-CONJelse if $\exists x, y : \text{domain}(root(G_x), x)$ and $: \text{domain}(root(G_{\bar{x}}, y) \text{ and } : \text{op*}(\text{``and''}, x) \in G_c$ and 32: $:op*(``and", y) \in G_c$ then # Comment: using connectives 'and' to connect two premises type = FRAME-CONJ33: else if $G_x \subset G_c$ then 34: 35: $d \leftarrow G_c - G_x$ if root(d) is a noun then 36: type = ARG-INS # Comment: inserting an argument. 37: 38: else type = FRAME-INS # Comment: inserting a phase (also annotated as ARG-INS). 39: end if 40: ---- ARG/PRED-GEN and Others -----41: # else if \exists :domain($root(G_c), y$) and ($root(G_c) \in G_x$ and $y \in G_{\bar{x}}$) or ($root(G_c) \in G_{\bar{x}}$ and $y \in G_x$) 42: then type = ARG/PRED-GEN 43: else 44: type = UNK45: end if 46: 47: end if

Prompts for automatic evaluation

Consistency:

You are a scoring expert in natural language reasoning. Given two premises and a conclusion, your goal is to evaluate whether the conclusion violates the premises. During your inference process, please only consider the information from the premises.

you can directly give your score (0 or 1) based on the following criteria:

0: the conclusion violates the premises.

1: the conclusion doesn't violate the premises.

The output format is just the score. You don't need to analyse the reasoning process.

Alignment:

You are a scoring expert. Given two premises, a conclusion, and an inference type, your goal is to evaluate whether the (premises, conclusion) pair is aligned with the inference type.

The following is the description of 10 inference types:

1. ARG-SUB: the conclusion is obtained by replacing one argument with another argument.

2. PRED-SUB: the conclusion is obtained by replacing one verb with another verb.

3. FRAME-SUB: the conclusion is obtained by replacing a frame of one of the premises with one from the other premise.

4. COND-FRAM: the conclusion is obtained according to the conditional premise with keyword "if".

5. ARG-INS: the conclusion is obtained by connecting an argument from one of the premises to a frame of the other.

6. FRAME-CONJ: the conclusion is obtained by using connectives to connect two premises.

7. ARG/PRED-GEN: a new ":domain" relation frame is created in the conclusion if both premise graphs differ by a single predicate/argument term.

8. ARG-SUB-PROP: one of the premises describes a "is made of" relationship between the entity in the other premise and its replacement.

9. IFT: the conclusion should be a conditional sentence.

10. EXAMPLE: the conclusion should contain the keyword "example".

When evaluating, some premises might not be able to deduce more than one conclusions. You can ignore those cases.

Finally, you can directly give your score (0 or 1) based on the following criteria:

0: the (premises, conclusion) pair is not aligned with the inference type.

1: the (premises, conclusion) pair is aligned with the inference type.

The output format is just the score. You don't need to analyse the reasoning process.

Table 17: Empirically designed prompt for automatically evaluating the controllability in Section 5.2.