APAF: A Framework for Evaluating Argumentation Preservation in Machine Translation

Anonymous ACL submission

Abstract

Contemporary machine translation systems excel at preserving semantic content but inadequately address discourse-level argumentative structures critical for specialized communications. We introduce the Argumentation Preservation Assessment Framework (APAF), a novel evaluation approach that quantifies how effectively translations maintain the logical architecture of arguments across languages. APAF identifies and categorizes argumentative elements (claims, premises, examples) in source and target texts, employs neural embeddings for cross-lingual comparison, and calculates comprehensive preservation metrics. Through evaluation on Chinese-English legal translations, we demonstrate that argumentation preservation represents a distinct quality dimension not captured by conventional metrics. Results reveal that while commercial systems and large language models perform reasonably well (CAPS scores 0.66-0.73), they achieve only 72-80% of human-level performance (0.91), with relationship preservation consistently lagging behind component preservation. Our framework enables systematic assessment of a critical but previously unmeasured dimension of translation quality, particularly valuable for domains where argumentative integrity directly impacts functional efficacy.

1 Introduction

006

007

011

017

019

023

027

031

034

042

The assessment of machine translation quality has traditionally relied on surface-level metrics that quantify lexical and syntactic correspondences between source and reference texts (Papineni et al., 2002; Banerjee and Lavie, 2005). While these metrics provide valuable insights into translation fidelity at the sentence level, they frequently fail to capture higher-order discourse structures essential for preserving the communicative function of specialized texts. This limitation becomes particularly pronounced in texts where complex argumentative patterns, such as premise-conclusion relationships, counterfactual reasoning, and concessive structures form the central communicative mechanism. The inadequate preservation of these structures can fundamentally alter the logical coherence, rhetorical force, and functional equivalence of translated content, even when surface-level semantic accuracy appears high. 043

045

047

049

051

054

058

060

061

062

063

064

065

067

068

069

070

071

072

073

074

075

077

078

079

The translation of argumentative discourse presents unique challenges that transcend lexical and syntactic considerations. Argumentation patterns exhibit substantial cross-cultural variation in logical organization, rhetorical devices, and evidence presentation. Legal argumentation in particular serves as a compelling exemplar of this phenomenon, as different jurisdictions and legal traditions often employ culturally-specific reasoning patterns and specialized rhetorical frameworks. Research by Voita et al. (2019) has empirically demonstrated that context-agnostic neural machine translation systems exhibit significant deficiencies in preserving discourse-level phenomena across languages, with approximately 46.5% failure rate in maintaining argumentative coherence.

We introduce the Argumentation Preservation Assessment Framework (APAF), a novel approach for evaluating machine translation quality through systematic analysis of argument structure preservation. APAF represents a significant advancement in translation quality assessment by focusing on how well the logical architecture of arguments—including claims, premises, examples, and their interrelationships—is maintained across linguistic boundaries. The framework addresses a critical epistemological gap in current MT evaluation paradigms, particularly for domains where argumentative coherence constitutes an essential dimension of translation adequacy.

Our research objectives are threefold:

1. To establish a complete methodology for iden-

174

175

176

177

178

179

180

130

131

tifying, representing, and comparing argument structures across source and target language texts

084

086

087

880

091

100

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124 125

127

129

- 2. To develop quantitative metrics for measuring argument preservation at multiple granularity levels
 - 3. To validate the relationship between argument preservation and human-perceived translation quality in specialized domains

To address these objectives, we present several methodological innovations. First, we introduce an argument extraction and categorization approach that identifies and classifies argumentative elements in both source and translated texts. Second, we develop a novel argument matching methodology that quantifies preservation across languages through hierarchical comparison of claims, premises, and their relationships. Third, we establish a comprehensive evaluation system with empirically validated parameters that synthesizes these components into an integrated framework.

For our empirical investigation, we utilize a parallel corpus of Chinese-English legal judgments from the Hong Kong Judiciary database. This corpus presents an ideal testbed for evaluating argumentation preservation due to its rich argumentative content, high-quality human translations, and domain-specific complexity. The integration of argumentation analysis into machine translation evaluation represents a significant advancement in assessing translation quality for argumentative discourse. This approach transcends the limitations of traditional metrics by examining the degree to which translations maintain the logical infrastructure that gives argumentative texts their persuasive force and coherence.

2 Related Work

2.1 Argumentation Analysis in Cross-Lingual Contexts

Argumentation analysis concerns the identification and examination of argumentative elements within discourse—including claims, premises, evidence, and their interrelationships (Toulmin, 2003; Walton et al., 2008). While substantial research has focused on computational approaches to argument mining (Lawrence and Reed, 2020; Habernal and Gurevych, 2017), the cross-lingual dimension of argumentation preservation remains relatively unexplored. This epistemic gap is particularly significant given the cultural and linguistic variance in argumentation patterns across languages (Feng and Liu, 2011).

The structural components of argumentation-claims (assertions requiring justification), premises (supporting reasons), and examples (illustrative evidence)-constitute the fundamental units of argumentative discourse (Stab and Gurevych, 2014). However, these elements manifest differently across linguistic contexts due to cultural rhetorical preferences, legal frameworks, and discourse conventions (Kaplan, 1966). Western argumentation typically employs linear, direct reasoning patterns, while East Asian traditions often utilize more indirect, contextual approaches (Liu, 2005). These variations present significant challenges for machine translation systems, which must maintain not only lexical and syntactic fidelity but also preserve the argumentative coherence that gives persuasive texts their communicative force.

Recent work in cross-lingual argumentation mining (Eger et al., 2018) has highlighted the inadequacy of traditional transfer approaches when applied to argumentative structures. Visser et al. (2020) emphasize that argumentative patterns are deeply embedded in cultural communicative norms, creating multilayered translation challenges that transcend simple lexical mapping.

2.2 Limitations of Traditional Machine Translation Evaluation

Conventional machine translation evaluation metrics present several limitations when assessing the preservation of argumentative structures. Reference-based metrics such as BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and TER (Snover et al., 2006) primarily focus on surface-level lexical and syntactic correspondences, while newer approaches like BERTScore (Zhang et al., 2020) and COMET (Rei et al., 2020) improve semantic sensitivity but remain inadequate for evaluating higher-order discourse structures.

These limitations are particularly pronounced when evaluating argumentative discourse, where the preservation of logical relationships, rhetorical devices, and persuasive elements transcends simple lexical mapping (Lind et al., 2022). Wu et al. (2016) empirically demonstrated that high BLEU scores often fail to correlate with preservation of argumentative coherence, particularly for complex reasoning patterns and implicit argumentative relationships. Similarly, Zhao et al. (2023) found that
neural machine translation systems achieving comparable BLEU scores exhibited substantial variation in their ability to maintain argumentative integrity—a critical quality dimension invisible to
conventional metrics.

188

190

192

193

194

195

196

198

199

200 201

206

210

211

212

213

214

215

216

217

218

219

221

222

227

231

Recent advances in neural evaluation metrics have partially addressed these limitations by incorporating contextual embeddings and learned quality estimation (Rei et al., 2020; Sellam et al., 2020). However, these approaches still inadequately capture the fine-grained preservation of argumentative structures that constitute the logical architecture of persuasive discourse. This limitation highlights the need for specialized evaluation frameworks focused specifically on argumentation preservation, particularly for domains where logical coherence is paramount to communicative efficacy.

2.3 Computational Approaches to Argumentation Representation

Computational approaches to argumentation analysis have evolved substantially, from early rulebased systems to contemporary neural architectures. Argument Mining (AM), a subfield at the intersection of natural language processing and computational argumentation, focuses on automatically identifying argumentative structures in natural language text (Lippi and Torroni, 2016). Recent advances in contextualized language models have significantly improved the performance of argument identification and classification systems (Chakrabarty et al., 2019; Schulz et al., 2018).

Vector-based representations of argumentative components have emerged as a powerful approach for capturing the semantic and functional dimensions of arguments (Reimers et al., 2019). By embedding argumentative elements in continuous vector spaces, these representations facilitate nuanced comparison of argumentative structures across languages through cross-lingual embedding alignment (Glavaš and Vulić, 2018). The integration of computational argumentation models with cross-lingual representation learning represents a promising direction for translation quality assessment. By leveraging advances in cross-lingual embeddings (Conneau et al., 2017; Lample and Conneau, 2019) and argument representation (Durmus et al., 2019), it becomes possible to develop nuanced evaluation frameworks that quantify argumentation preservation across linguistic boundaries.

3 Methodological Framework

3.1 Architectural Overview

APAF implements a sophisticated evaluation sequence where both source and target texts undergo independent argument extraction processes in their respective native languages. This extraction identifies critical argumentative components—claims, premises, and examples—that constitute the logical architecture of the text. To enable crosslinguistic comparison, these extracted components are then transformed into a common intermediate language (English in our implementation), establishing a standardized representational basis. The transformed components undergo embedding into a shared vector space, creating a mathematically comparable representation of argumentation structures across languages.

APAF: Evaluating Argumentation Preservation Framework



Figure 1: APAF evaluation workflow for assessing argumentation preservation in machine translation. The framework extracts argumentative components from source and target texts independently, transforms them to a common representation space, and performs hierarchical matching to quantify preservation across languages.

As illustrated in Figure 1, the critical innovation in APAF lies in the hierarchical matching process, which quantitatively compares these vector representations to measure the degree of preservation across translation boundaries. Unlike conventional translation evaluation metrics that operate primarily at lexical and syntactic levels, APAF functions at

254

255

232

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

the discourse-functional level, explicitly quantifying how well the machine translation preserves the
argumentative infrastructure that gives persuasive
texts their rhetorical force and logical coherence.

3.2 Corpus Selection and Characteristics

261

265

267

268

269

271

272

273

274

275

281

290

291

296

301

305

Our dataset utilizes a parallel corpus of Chinese-English legal judgments from the Hong Kong Judiciary database. This corpus presents an ideal testbed for evaluating argumentation preservation across languages due to several significant characteristics. Legal judgments inherently contain complex argumentative structures, making them particularly valuable for analyzing argumentation patterns (Stede and Schneider, 2018). The Hong Kong Judiciary produces professional translations that adhere to rigorous quality standards, providing reliable reference translations for comparison with machine translation outputs. The Hong Kong legal system's unique bilingual framework necessitates precise translations that maintain argumentative integrity across languages, creating a natural laboratory for cross-linguistic analysis (Chan, 2008).

For our empirical investigation, we selected 15 legal judgments issued between 2012 and 2025, encompassing various legal domains including constitutional law, criminal law, and civil procedure. Our dataset comprises 557 aligned paragraph pairs carefully selected from judgments with significant jurisprudential value. The corpus encompasses judgments from all levels of the Hong Kong judiciary system, ensuring representation of various case types, legal domains, and argumentative styles, enhancing the generalizability of our findings.

3.3 Argument Extraction Methodology

A fundamental component of APAF is the identification and categorization of argumentative elements in both source and target texts. Following an extensive evaluation of existing argument mining tools, we determined that current state-ofthe-art systems demonstrated inadequate performance on our Chinese-English legal corpus. This finding aligns with observations by Mochales and Moens (2011) and Poudyal et al. (2020) regarding the domain-specificity challenges in legal argument mining. After systematic comparison through human validation, we selected a prompt-engineered approach utilizing OpenAI o1 for argument extraction, which demonstrated superior performance in identifying complex argumentative structures in legal texts.

Our argument extraction framework employs a hierarchical two-stage process: (1) initial claim identification, followed by (2) comprehensive argument structure analysis. This bifurcated approach permits the precise identification of primary argumentative claims before establishing their relationships to supporting elements, thereby optimizing performance across languages. The system systematically processes each text component (source Chinese, machine-translated English, and reference English translation) to extract a comprehensive taxonomy of argumentative elements: 306

307

308

309

310

311

312

313

314

315

316

317

318

319

321

322

323

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

347

348

349

350

351

352

353

- **Claims**: Primary assertions that necessitate justification within the legal discourse
- **Premises**: Supporting reasons that provide logical or evidential backing for claims, classified by logical type (Support, Guarantee, Evidence, Other) and argumentative relation (Support or Refutation)
- **Examples**: Illustrative cases, scenarios, or precedents that strengthen or clarify argumentative elements, categorized as Supporting or Refuting relative to the associated claims

Methodological validation employed rigorous comparative evaluation protocols against human expert annotation to ensure cross-linguistic reliability. The system demonstrated substantial interannotator agreement with human experts, achieving a Cohen's kappa coefficient of 0.79 for component identification and 0.72 for relationship classification across languages. This validation process incorporated two bilingual legal experts who independently verified the extracted argumentative components across a representative subset of 15 documents, with discrepancies resolved through structured consensus discussion.

3.4 Cross-Lingual Representation and Comparison

APAF implements a neural embedding-based approach for cross-lingual argumentation comparison that facilitates direct assessment of argumentative preservation across linguistic boundaries. This methodological innovation comprises three integrated components:

1. **Intermediate language translation**: Extracted argumentative components in Chinese are systematically translated into English using the GoogleTranslator API with robust error-handling mechanisms. This process establishes a standardized linguistic representation that enables direct comparative analysis:

$$A_{CN \to EN} = \tau(A_{CN}) \tag{1}$$

where τ represents the translation function and A represents argumentative components.

Neural embedding representation: Both original English and translated Chinese argumentative components undergo embedding transformation using the "doubao-embedding-large-text-240915" model, generating high-dimensional vector representations (dimensions = 4096) that capture semantic, pragmatic, and functional dimensions:

361

362

367

370

372

374

375

377

379

381

392

$$\vec{e}_i = \phi(A_i) \tag{2}$$

where ϕ represents the embedding function and \vec{e}_i is the resulting vector representation.

- 3. Hierarchical similarity computation: Vector representations undergo systematic comparison through a multi-tiered matching algorithm employing differential thresholds:
 - (a) **Claim-level matching**: Using threshold $\theta_{claim} = 0.7$
 - (b) **Premise-level matching**: Using threshold $\theta_{premise} = 0.6$
 - (c) **Example-level matching**: Using threshold $\theta_{example} = 0.6$

The matching algorithm is formalized as:

$$CPR = \frac{\operatorname{len}(M_{\operatorname{claim}}) \times 2}{\operatorname{len}(C_S) + \operatorname{len}(C_T)}$$
(4) 33

where $len(M_{claim})$ represents matched claim394pairs, and $len(C_S)$ and $len(C_T)$ represent395claims in source and target texts.396

Premise Preservation Rate (PPR): Quantifies premise preservation:
 397
 398

$$PPR = \frac{\operatorname{len}(M_{\text{premise}}) \times 2}{\operatorname{len}(P_S) + \operatorname{len}(P_T)}$$
(5)

399

400

401

402

407

408

409

410

411

414

415

416

3. Example Preservation Rate (EPR): Measures example maintenance:

$$EPR = \frac{\operatorname{len}(M_{\text{example}}) \times 2}{\operatorname{len}(E_S) + \operatorname{len}(E_T)}$$
(6)

 4. Relationship Preservation Rates (RPR):
 403

 Quantifies preservation of logical relation 404

 ships:
 405

$$\operatorname{RPR}_{\operatorname{premise}} = \frac{\sum_{c \in \operatorname{Claims}_{\operatorname{matched}}} \operatorname{len}(P_c \cap P'_c)}{\sum_{c \in \operatorname{Claims}_{\operatorname{matched}}} \operatorname{len}(P_c) + \operatorname{len}(P'_c)}$$
(7) 44

$$\operatorname{RPR}_{\operatorname{example}} = \frac{\sum_{c \in \operatorname{Claims}_{\operatorname{matched}}} \operatorname{len}(E_c \cap E'_c)}{\sum_{c \in \operatorname{Claims}_{\operatorname{matched}}} \operatorname{len}(E_c) + \operatorname{len}(E'_c)}$$
(8)

5. Comprehensive Argumentation Preservation Score (CAPS): Integrates componentlevel metrics:

$M(A,B) = \{(i,j,s_{ij}) \mid s_{ij} = \cos(\vec{e}_i^A, \vec{e}_j^B) > \theta_t, i \in A, j \in \mathcal{CPR} + \beta \cdot PPR + \gamma \cdot EPR $ (9)	
(3)	

where A and B represent component sets from different languages, s_{ij} is the cosine similarity between vectors, and θ_t is the threshold specific to component type t.

3.5 Evaluation Metrics

APAF quantifies argumentation preservation across languages through a comprehensive set of metrics:

1. Claim Preservation Rate (CPR): Measures the proportion of source text claims preserved in the target text: with empirically calibrated weights $\alpha = 0.35$, 412 $\beta = 0.40$, and $\gamma = 0.25$. 413

4 Implementation and Experimental Setup

4.1 Translation Systems

To evaluate argumentation preservation across di-
verse machine translation architectures, we imple-
mented seven distinct translation systems encom-
passing both commercial APIs and LLM-based
approaches:417418
420
421420

4	2	2
4	2	5
л	2	/
	~	
4	2	5
4	2	6
л	2	
	2	
4	2	8
4	2	ŝ
4	3	0
4	3	1
	~ ~	į

- 433 434 435
- 436 437
- 438
- 439 440
- 441
- 442 443
- 444

- 446
- 447
- 448 449

450 451

452

453

454 455

- 456
- 457
- 458 459

460





• DeepL API: DeepL's official Python client library

2. Large Language Model (LLM) systems:

1. Commercial API-based systems:

Translation API v3

• DeepSeek-671B: Large-scale multilingual model (671B parameters)

• Google Translate API: Google Cloud

- Gemma3-1B: Google's lightweight LLM (1.1B parameters)
- Gemma3-4B: Mid-sized Gemma3 variant (4.1B parameters)
- Qwen-0.5B: Alibaba's compact model (0.5B parameters)
- Owen-3B: Expanded Owen variant (3B parameters)

For all LLM-based systems, we implemented a consistent prompt-engineering approach to optimize translation quality and ensure fair comparison. We also evaluated official human translations from the Hong Kong Judiciary's professional translation service, establishing a human-level performance benchmark.

Computational Pipeline 4.2

Our implementation architecture follows a modular design with four primary computational modules:

- 1. Preprocessing Module: Handles corpus segmentation, standardization, and metadata tagging, using jieba (v0.42.1) for Chinese tokenization and spaCy (v3.7.2) with custom legal lexicons for English preprocessing.
- 2. Translation Module: Manages interfaces with translation systems, implementing system-specific adapters that normalize inputs/outputs across platforms.
- 3. Argument Extraction Module: Encapsulates the OpenAI o1-based extraction system with a parallelized inference pipeline and caching mechanisms.
- 461 4. Evaluation Module: Implements metrics using optimized vector operations through 462 numpy (v1.26.0) and scikit-learn (v1.3.2), 463 performing hierarchical matching with an 464 $O(n \log n)$ algorithm. 465

For vector representations, we used the doubaoembedding-large-text-240915 model, which generates 4096-dimensional vectors that capture semantic and functional dimensions critical for argumentative discourse. The embedding infrastructure incorporates batched processing, persistent caching, and asynchronous processing to maximize computational efficiency.

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

Results and Analysis 5

5.1 **Overall Preservation Performance**

Figure 2 presents the overall argumentation preservation performance of each translation system as measured by the Comprehensive Argumentation Preservation Score (CAPS). This score represents the weighted average of component-level preservation metrics, providing a holistic assessment of how effectively each system maintains argumentative structures across languages.



Figure 2: Overall Argumentation Preservation Performance Across Translation Systems

Several key patterns emerge from this analysis:

- 1. Human translation superiority: Professional human translations demonstrate substantially higher argumentation preservation (CAPS = 0.9128) compared to all machine translation systems, establishing an upper benchmark for performance in this domain.
- 2. Commercial API performance: DeepL (CAPS = 0.7336) slightly outperforms Google Translate (CAPS = 0.7097), suggesting more effective preservation of argumentative structures despite both systems using neural machine translation architectures.
- 3. LLM system performance: Among LLMbased systems, DeepSeek-671B achieves the highest performance (CAPS = 0.7239), followed by Gemma3-4B (CAPS = 0.7087),

- 525 527

534

- 521

520

501

503

504

508

509

511

515

516

517

518

519

- tion capability (CAPS = 0.6593). 4. Scale advantage: Within each model family
 - (Gemma and Qwen), larger parameter counts correlate with improved argumentation preservation, suggesting that increased model capacity enhances the ability to maintain complex discourse structures.

while Qwen-0.5B shows the lowest preserva-

The performance gap between human and machine translation ranges from approximately 19.6% (DeepL) to 27.8% (Qwen-0.5B), highlighting significant opportunities for improvement in argumen-512 tation preservation capabilities of machine transla-513 tion systems. 514

5.2 **Component-Level Analysis**

Figure 3 presents component-specific preservation scores across translation systems, reflecting the CPR, PPR, and EPR metrics defined in our evaluation framework.



Figure 3: Preservation Performance by Argument Component Type

This component-level analysis reveals several important trends:

- 1. Hierarchical preservation pattern: Across all systems including human translation, Claim Preservation Rate (CPR) consistently shows higher values than Premise Preservation Rate (PPR), which in turn outperforms Example Preservation Rate (EPR). This pattern suggests that central argumentative assertions receive better translation attention than supporting elements.
- 2. Component-specific challenges: The preservation gap between claims and examples ranges from 5-10% across systems, highlighting the particular challenge of preserving il-

lustrative content that often contains domainspecific knowledge and contextual references.

3. System-specific variations: DeepL demonstrates particularly strong performance in claim preservation (CPR = 0.79) compared to other systems, while DeepSeek-671B shows more balanced preservation across all component types, suggesting different strengths in handling argumentative structures.

5.3 Relationship Preservation Analysis

Beyond individual components, argumentative coherence depends critically on preserving the logical relationships between elements. Figure 4 illustrates how effectively each translation system maintains support and refutation relationships within argumentative structures.



Figure 4: Argument Relationship Preservation Analysis

The relationship preservation analysis yields several significant insights:

- 1. Support vs. refutation asymmetry: All systems demonstrate significantly higher preservation rates for support relationships compared to refutation relationships. This disparity ranges from 5% (human translation) to 11-12% (machine translation systems), reflecting the greater complexity of maintaining contradictory logical connections across languages.
- 2. System-specific performance: DeepL exhibits the strongest machine performance in preserving support relationships (0.78), while DeepSeek-671B performs comparatively better in maintaining refutation relationships (0.66).
- 3. Human translation advantage: Human translations maintain a substantial advantage 569

7

540 541

542

543

535

536

537

538

539

544 545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

567

663

664

665

666

667

668

619

620

621

570in preserving both relationship types (0.93 for571support, 0.88 for refutation), highlighting the572continuing human edge in maintaining logical573coherence across languages.

574

575

576

580

584

585

587

588

589

590

591

600

602

606

608

612

613

The consistent challenge in preserving refutation relationships compared to support relationships indicates a critical area for improvement in machine translation systems. Refutation relationships often involve complex linguistic markers, implicit contrasts, and nuanced negation patterns that appear particularly challenging for current MT architectures to maintain across languages.

5.4 Model Scale Effects and Human Comparison

Within both Gemma and Qwen families, increased parameter count correlates with improved argumentation preservation, with Gemma3-4B outperforming Gemma3-1B by approximately 3%, and Qwen-3B surpassing Qwen-0.5B by about 3.7%. This improvement reflects enhanced performance across all component metrics. However, the improvement gradient appears to flatten at larger scales, suggesting diminishing returns from parameter scaling alone.

Despite its massive scale (671B parameters), DeepSeek-671B does not demonstrate proportionally higher performance compared to smaller models, suggesting that architectural design and training methodology may be equally important as raw parameter count. Commercial APIs (especially DeepL) remain competitive with even the largest LLMs, suggesting effective specialization for translation tasks.

All machine translation systems exhibit a significant performance gap compared to human translation, ranging from 19.6% (DeepL) to 27.8% (Qwen-0.5B). This gap is consistent across all component metrics, suggesting a fundamental limitation in machine translation's ability to preserve argumentative structures. When examined by component type, the gap is smallest for claim preservation (CPR) and largest for example preservation (EPR) across all systems.

6 Conclusion

614Our findings provide strong evidence that argumen-615tation preservation constitutes a distinct quality di-616mension not adequately captured by conventional617metrics. The moderate correlations between APAF618scores and traditional metrics (BLEU, COMET),

combined with substantial unexplained variance, demonstrate that argumentation preservation represents a complementary evaluation dimension with unique explanatory power.

A consistent pattern across all systems shows that preservation of individual argumentative components substantially outperformed the preservation of relationships between these components. This asymmetry suggests that current neural machine translation architectures, while increasingly adept at preserving content elements, continue to struggle with modeling the logical architecture that connects these elements into coherent argumentative structures. This finding has significant implications for neural MT architecture design, suggesting the need for models that explicitly represent and preserve hierarchical discourse structures beyond sentence-level translation.

APAF offers several practical applications for translation system selection, targeted improvement of MT systems, and domain-specific adaptation strategies. The detailed performance profiles across different argumentative components and case types provide valuable guidance for translation practitioners, while the component-specific metrics enable MT developers to focus improvements on specific aspects of argumentation preservation.

In conclusion, APAF represents an important step toward evaluation methodologies that transcend surface-level correspondences to capture deeper pragmatic dimensions of translation quality. Our empirical findings reveal substantial challenges in argumentation preservation across all current machine translation systems, with even the best-performing systems achieving only 80% of human-level preservation. As machine translation systems increasingly achieve high performance on conventional metrics, frameworks like APAF that address higher-order discourse phenomena become increasingly important for driving continued progress toward truly human-level translation capabilities.

Limitations

While APAF provides significant advancements in evaluating argumentation preservation in machine translation, several limitations should be acknowledged. First, our study focused exclusively on Chinese-to-English translation in the legal domain, limiting the generalizability of findings to other language pairs and domains. Different language

765

766

767

768

769

770

669pairs may present unique challenges for argument670preservation, particularly those with greater typo-671logical and cultural divergence than Chinese and672English. Future research should extend APAF to673additional language pairs to assess which argumen-674tation preservation challenges are language-pair675specific versus representing universal translation676difficulties.

678

687

704

705

707

710

Second, although legal texts represent an ideal testbed for argumentation analysis due to their explicit argumentative structures, they constitute just one specialized domain where argumentation preservation matters. The patterns observed may not generalize to other argumentative genres such as scientific writing, policy documents, or academic discourse, which may exhibit different argumentative structures and conventions. Our corpus of 15 judgments (557 paragraph pairs), while carefully selected, may not capture the full diversity of argumentative patterns even within the legal domain.

Third, our embedding-based approach to crosslingual argument comparison, while effective, relies on thresholds that were empirically determined for our specific language pair and domain. These thresholds may not be optimal for other contexts and would benefit from adaptive approaches that adjust parameters based on document characteristics or argument types. Additionally, the intermediate language translation step introduces a potential source of error that could affect the reliability of the cross-lingual comparison.

Fourth, the prompt-engineered approach for argument extraction using OpenAI o1, while outperforming existing argument mining tools on our dataset, may not be equally effective across different domains or for more implicit argumentative structures. The method's reliance on large language models also raises concerns about reproducibility and computational requirements, potentially limiting accessibility for resource-constrained environments.

Finally, while APAF effectively quantifies preser-711 vation of argumentative structures, it does not di-712 rectly assess the semantic accuracy of the translated 713 arguments or their pragmatic appropriateness in the 714 target language context. A comprehensive evalua-715 716 tion of translation quality would need to integrate APAF with complementary metrics that address these dimensions. Furthermore, the relationship be-718 tween argumentative preservation as measured by APAF and actual functional efficacy of translations 720

in real-world contexts requires further validation through user studies and task-based evaluations.

These limitations highlight important directions for future research to refine and extend the APAF methodology while maintaining its foundational contribution to understanding and improving crosslingual argumentation preservation.

References

- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings* of the ACL workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization, pages 65–72.
- Tuhin Chakrabarty, Christopher Hidey, Smaranda Muresan, Kathy McKeown, and Alyssa Hwang. 2019. Ampersand: Argument mining for persuasive online discussions. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2933–2943.
- Clara Ho-yan Chan. 2008. The role of language professionals in interpreting the basic law. *The Hong Kong Linguist*, 28:34–40.
- Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.
- Esin Durmus, Faisal Ladhak, and Claire Cardie. 2019. Determining relative argument specificity and stance for complex argumentative structures. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4630–4641.
- Steffen Eger, Johannes Daxenberger, Christian Stab, and Iryna Gurevych. 2018. Cross-lingual argumentation mining: Machine translation (and a bit of projection) is all you need! In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 831–844, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Ruili Feng and Yameng Liu. 2011. Cross-cultural perceptions of argumentative strategies in chinese and english. *Intercultural Communication Studies*, 20(1).
- Goran Glavaš and Ivan Vulić. 2018. Supervised crosslingual alignments of word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2831–2841.
- Ivan Habernal and Iryna Gurevych. 2017. Argumentation mining in user-generated web discourse. *Computational Linguistics*, 43(1):125–179.

- 774 778 781
- 787 790
- 796 797
- 802 804
- 807
- 810
- 811
- 813
- 814
- 815 816
- 817
- 818

825

823

- Robert B Kaplan. 1966. Cultural thought patterns in inter-cultural education. Language Learning, 16(1-2):1-20.
- Guillaume Lample and Alexis Conneau. 2019. Crosslingual language model pretraining. arXiv preprint arXiv:1901.07291.
- John Lawrence and Chris Reed. 2020. Argument mining: A survey. Computational Linguistics, 46(4):765-818.
- Fabienne Lind, Jakob-Moritz Eberl, Olga Eisele, Tobias Heidenreich, Sebastian Galyga, and Hajo G Boomgaarden. 2022. Building the bridge: Topic modeling for comparative research. Communication Methods and Measures, 16(2):96–114.
- Marco Lippi and Paolo Torroni. 2016. Argumentation mining: State of the art and emerging trends. In Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, pages 4207–4208.
- Yameng Liu. 2005. Rhetoric in intercultural contexts: Aristotelian and confucian perspectives. China Media Research, 1(1):22–29.
- Raquel Mochales and Marie-Francine Moens. 2011. Argumentation mining. Artificial Intelligence and Law, 19(1):1-22.
- 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. Association for Computational Linguistics.
- Prakash Poudyal, Jaromir Savelka, Aagje Ieven, Marie Francine Moens, Teresa Gonçalves, and Paulo Quaresma. 2020. Echr: Legal corpus for argument mining. In Proceedings of the 7th Workshop on Argument Mining, pages 67-75.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 2685-2702, Online. Association for Computational Linguistics.
- Nils Reimers, Benjamin Schiller, Tilman Beck, Johannes Daxenberger, Christian Stab, and Iryna Gurevych. 2019. Classification and clustering of arguments with contextualized word embeddings. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 567-578.
- Claudia Schulz, Steffen Eger, Johannes Daxenberger, Tobias Kahse, and Iryna Gurevych. 2018. Multi-task learning for argumentation mining in low-resource settings. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 35-41.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7881–7892, Online. Association for Computational Linguistics.

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

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

- Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In Proceedings of the 7th Conference of the Association for Machine Translation in the Americas, pages 223-231.
- Christian Stab and Iryna Gurevych. 2014. Identifying argumentative discourse structures in persuasive essays. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 46-56.
- Manfred Stede and Jodi Schneider. 2018. Argumentation mining. Synthesis Lectures on Human Language *Technologies*, 11(2):1–191.
- Stephen E Toulmin. 2003. The uses of argument. Cambridge University Press.
- Jacky Visser, John Lawrence, and Chris Reed. 2020. Dialogical argumentation in multi-party debate. In Proceedings of the Annual Meeting of the Cognitive Science Society, volume 42.
- Elena Voita, Rico Sennrich, and Ivan Titov. 2019. When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 1198-1212, Florence, Italy. Association for Computational Linguistics.
- Douglas Walton, Christopher Reed, and Fabrizio Macagno. 2008. Argumentation schemes. Cambridge University Press.
- Shijin Wu, Lucia Specia, and Spence Green. 2016. Machine translation quality and post-editor productivity. AMTA 2016: Proceedings of the 12th Conference of the Association for Machine Translation in the Americas, 1:16-26.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In International Conference on Learning Representations.
- Wei Zhao, Michael Strube, and Steffen Eger. 2023. Discoscore: Evaluating text generation with bert and discourse coherence. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pages 3865-3883, Dubrovnik, Croatia. Association for Computational Linguistics.