

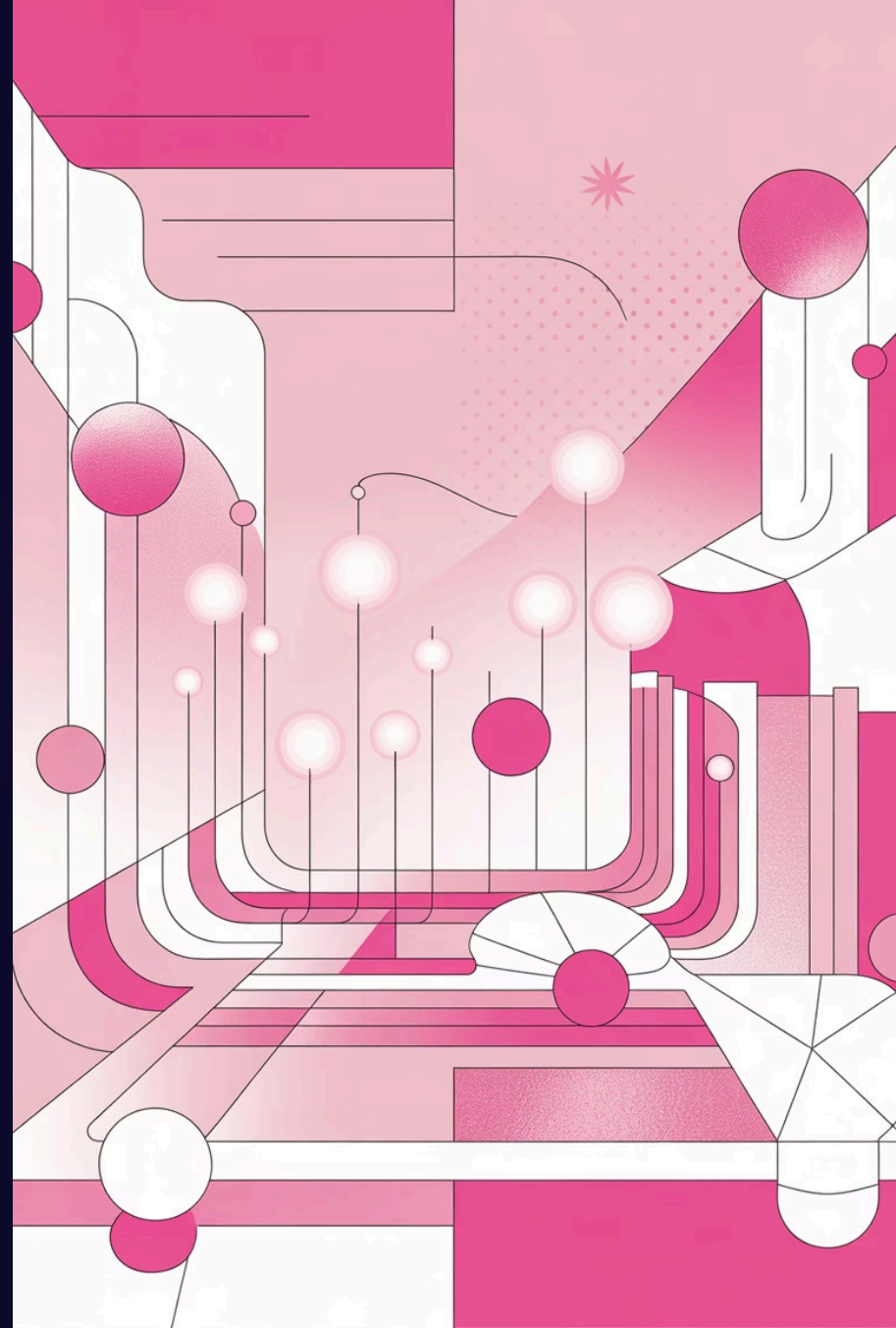
Natural Language Processing

Cross-lingual Semantic Textual Similarity

Presented by Amirreza Sharifzade

Cross-lingual Semantic Textual Similarity

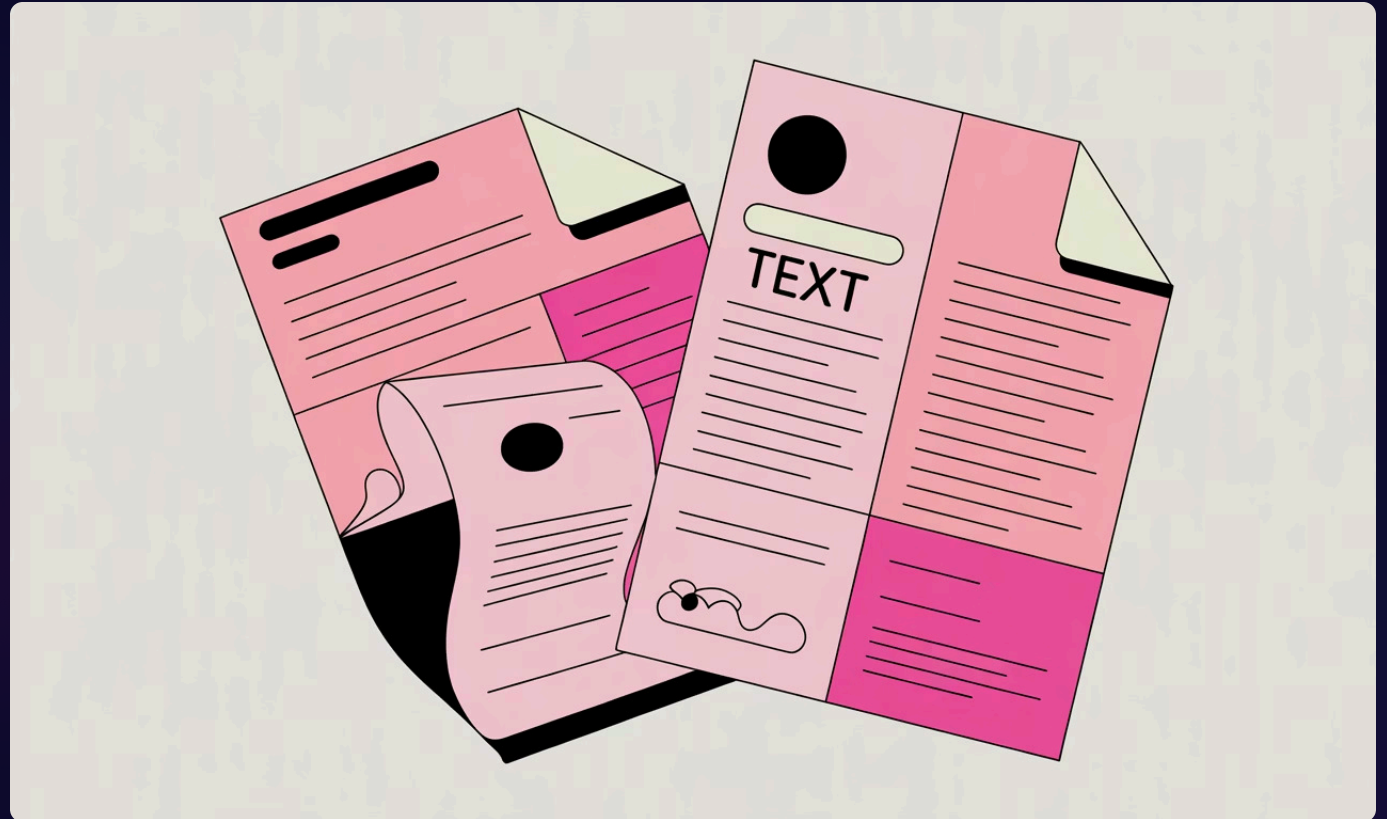
A fundamental task in multilingual NLP that measures semantic similarity between texts written in different languages, enabling language-independent understanding.



What Is Semantic Textual Similarity?

Core Concept

Semantic Textual Similarity (STS) estimates how closely two sentences match in meaning. The similarity is represented as a numerical score, focusing on semantic understanding rather than exact word overlap.



Meaning-Focused

Looks beyond surface-level word matching to capture deeper semantic relationships

Numerical Score

Produces a similarity value that quantifies the degree of semantic alignment

Building Block

Serves as foundation for numerous NLP applications and systems

Texts are written in distinct languages, requiring cross-lingual understanding

2

3

The key challenge: capturing shared meaning despite linguistic differences, cultural nuances, and vocabulary gaps between languages.

Why CL-STS Matters

Global Information Access

Enables multilingual search and information retrieval across language barriers

Reduced Translation Dependency

Operates without explicit machine translation, reducing error propagation

Fairness & Inclusivity

Improves NLP systems for speakers of diverse languages

Security Applications

Detects threats and misinformation across multiple languages simultaneously



Real-World Applications



Cross-lingual Search

Search engines that retrieve relevant documents regardless of query or document language



Plagiarism Detection

Identifies copied content across different languages in academic and publishing contexts



Content Moderation

Moderates harmful content in multilingual platforms by detecting similar problematic messages



Misinformation Detection

Tracks fake news and disinformation spreading across language boundaries



MT Evaluation

Assesses machine translation quality by comparing source and target sentence meaning



Social Media Analysis

Monitors coordinated campaigns and global discussions across diverse linguistic communities



Position in the NLP Landscape



Sentence Embeddings

Produces dense vector representations that capture semantic meaning



Transfer Learning

Enables knowledge transfer between languages and tasks



Information Retrieval

Power's multilingual search and document matching systems



Language Models

Connected to multilingual Transformer architectures

Modern Technical Approaches

LaBSE (2022)

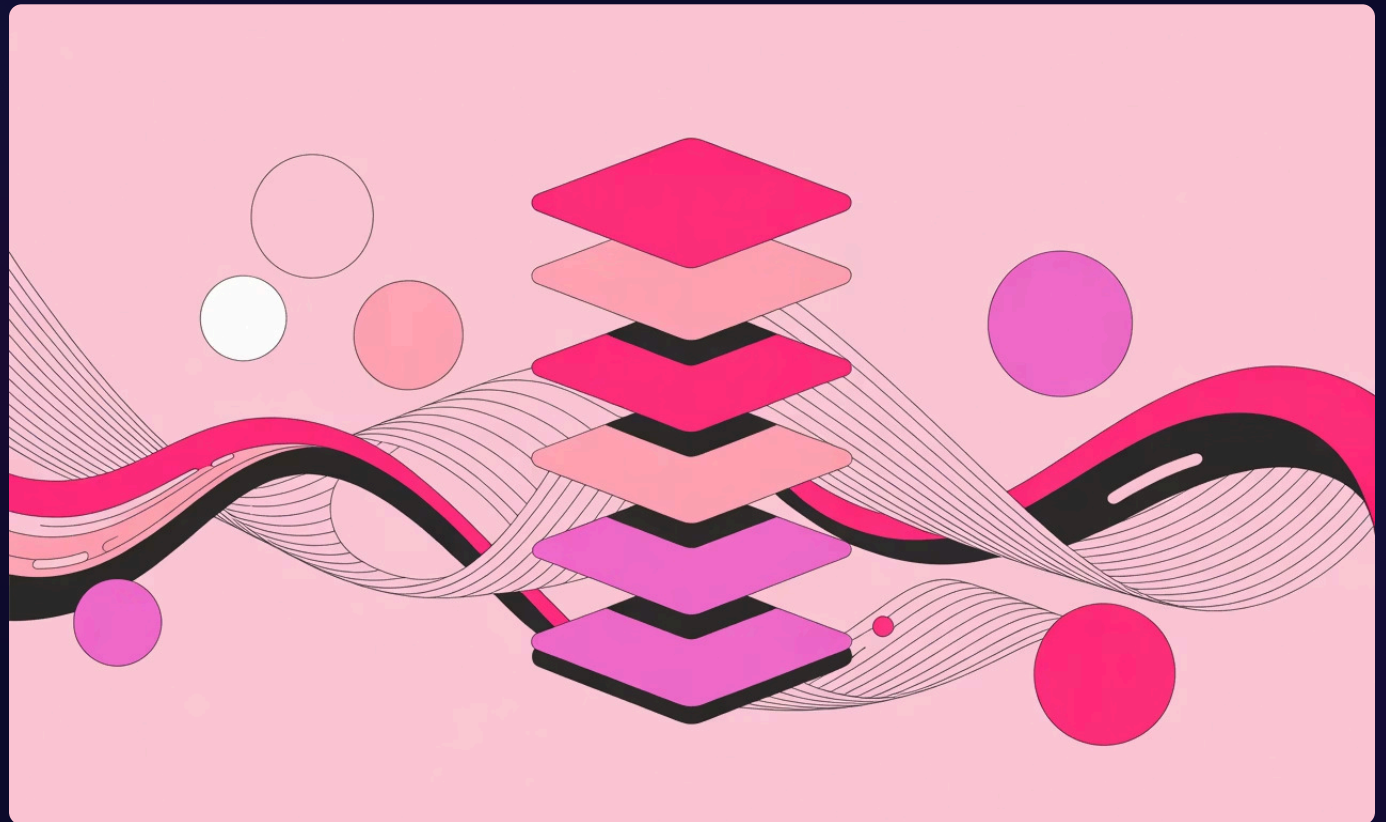
Language-agnostic BERT Sentence Embedding maps sentences from different languages into shared vector space. No translation needed at inference time. Strong performance on CL-STS benchmarks.

mSBERT (2022)

Multilingual extension of Sentence-BERT using knowledge distillation from English STS data. Effective for low-resource languages with performance drops for distant language pairs.

XLM-R Methods (2023)

Built on multilingual Transformer architectures pre-trained on massive multilingual corpora. Fine-tuned on CL-STS datasets, outperforming traditional cross-lingual methods.



📌 **Key insight:** Embedding-based methods using shared vector spaces dominate current research, replacing older translation-based and alignment-based approaches.



Research Challenges

1

Limited Labeled Data

Shortage of annotated STS datasets for many languages, especially low-resource ones

2

Cultural & Semantic Differences

Meaning doesn't always translate directly due to cultural context and linguistic structures

3

High-Resource Bias

Models perform better on widely-spoken languages, disadvantaging others

4

Computational Costs

Training multilingual models requires massive datasets and significant compute resources

Current Trends & Future Directions



Key Takeaways

Foundation Task

CL-STS is a fundamental building block in multilingual NLP

Embedding-Based

Shared sentence embeddings now dominate approaches

Transformer-Powered

Modern architectures significantly improved performance

Active Research

Low-resource languages remain key challenge

CL-STS removes language barriers in semantic understanding, making it essential for global AI systems with strong real-world impact.

