**Academic Education Publishing House**
-AEPH-

# Corpus-Aided Translation for Culture-Loaded Terms: Applications and Limitations in Audiovisual Subtitling

**Tianyu Gao**[*]

*School of Language and Communication, Beijing Technology and Business University, Beijing, China*
*Corresponding Author*

**Abstract: This study examines the advantages and limitations of corpus-aided translation in addressing Beijing-specific, culture-loaded terms within audiovisual subtitles. While corpus tools offer empirical support through authentic language patterns, concordance lines, and translation norms, their application in translating culturally nuanced expressions—such as Hutong, Siheyuan, and Douzhe—is often constrained by low frequency, data obsolescence, and lack of audiovisual context. Drawing on theoretical frameworks including Skopos theory, Pedersen's ECR strategies, and Descriptive Translation Studies, the paper critically assesses the role of corpus resources in balancing linguistic consistency with cultural fidelity. Empirical cases and bilingual subtitle corpora are analyzed to show how translation decisions are informed—or limited—by corpus evidence. The findings emphasize that corpus tools are practical foundations for subtitle translation. However, they must be complemented by human judgment, contextual awareness, and multimodal integration to preserve linguistic specificity and translation functionality. The research aims to develop dynamic, genre-specific, and multimodal-enriched subtitle corpora for better audiovisual translation (AVT).**

**Keywords: Corpus-Aided Translation; Audiovisual Translation; Subtitles; Culture-loaded Terms; Beijing Dialect**

## 1. Introduction

The purpose of the present study is to ascertain the efficacy and limits of corpus-assisted translation available for culturally loaded terms in AVT, namely those which originated from Beijing. While corpora can offer important patterns and precedents, if and how they could be applied for low-frequency culturally loaded expressions are questions open for argument. This is especially true with culturally specific material and dialectal variation that does not so often appear in corpora that are readily available. Such items are introduced too little with audiovisual context.

This paper analyzes whether and in which ways corpus tools may assist or otherwise provide obstacles to the translation into subtitles of Beijing-specific linguistic expressions. Past research and real subtitling practices will give an opportunity to examine the potential and limitations offered by corpus-aided methods. The methodology includes a critical literature review and an analytical discussion of representative examples.

## 2. Literature Review

### 2.1 Corpus-Aided Translation: an Evolving Paradigm

Translation studies is undergoing rapid development of corpus-based translation, particularly with the advent of computerized materials that are making it possible to conduct large-scale analyses of actual language use. According to Laviosa, corpora provide valuable resources for finding patterns, changes in translation, and normativity in bilingual writing [1]. Corpus usage makes it possible for translation decision-makers to gain knowledge of frequency, collocations, and genre-specific use.

In their 2024 work, Corbett and Li examined how bilingual subtitle corpora trigger pragmatic markers and colloquial tones for streaming content across languages [2]. Corpus tools have the potential to significantly promote the consistency of subtitle making and the naturalness of the production process when well designed and domain-specific. Although such tools can be utilized on a wide range of content types and cultural settings, their comprehensiveness and limitations remain disproportionate.

## 2.2 The Role of Corpora in Audiovisual Translation

The main difference of AVT from traditional translation is that AVT is carried out under technical, spatial, and temporal constraints. It is particularly important to make subtitles brief without losing meaning and cultural identity. Among the outcomes of the study of Pavesi is the identification of repetition of translation solutions by genre and the investigation of subtitle norms [3]. The multilingual TED Talks corpus and OpenSubtitles are useful for translators operating under such constraints. Corbett and Li's research into subtitle norms discovered that although corpora assist in the identification of default approaches—like omission, generalization, or substitution—these strategies hinder culturally embedded meaning [2]. This is particularly challenging for extra-linguistic cultural references (ECRs), where literal or neutral translation conveys inadequate sociocultural implications.

## 2.3 Extra-linguistic Cultural References and Corpus Use

ECRs are elements based on the social, historical, or geographical context of the culture and include traditions, institutions, food, and dialectical expressions. ECRs are difficult to translate in AVT since they carry implied meaning and rely on background cultural knowledge. Pedersen states that good ECR rendition is often a hybrid solution that combines corpus insights with cultural expertise [4].

Corpus resources provide hints at how certain ECRs were handled in subtitles previously, typically by comparable bilingual segments. However, such corpus resources are short of efficacy when dealing with underrepresented or low-frequency items.

## 2.4 Beijing-Specific Cultural Terms in Translation Research

Specific to Beijing, growing scholarly focus is being given on the emergence of terms because of the global spread of Chinese television dramas and films. Yu studied the subtitles of those popular dramas focusing on the dialect expressions from Beijing with the general conclusion that corpus base assessments often fail when it comes to words that hold in-depth cultural meaning [5]. For example, terms like douzhe, meaning usually a traditional opera lover or comic performer, might not have equivalents as such and would need context in order to understand. This paper sheds light that corpus data should be integrated along with ethnographic or contextual interpretations.

A recently done study by Yao et al. shows the limitations of national-level corpora toward regional cultural nuances [6]. In generalizations, it has been illustrated that large corpora, like Chinese National Corpus, provides loss of cultural specificity.

## 2.5 Summary of Findings

The reviewed literature establishes the foundational role of corpora in facilitating systematic, evidence-based translation strategies. In AVT, particularly in subtitles, corpora enhance linguistic accuracy and functional consistency. However, when it comes to their role in translating ECRs, and more specifically Beijing-specific linguistic terms, remains constrained by limitations in coverage, frequency, and context.

## 3. Theoretical Framework

### 3.1 Skopos Theory and Functionalist Approaches

According to skopos theory introduced by Vermeer, translation is governed by its purpose (skopos) [7]. Translation strategy is selected in this approach based on the role a target text is to play within its target culture. In audiovisual translation (AVT), and subtitles in particular, this entails foregrounding readability, cultural relevance, and audience expectation over strict fidelity to the source text. Handling Beijing-specific culture-loaded terminology, the Skopos principle encourages translators to balance authenticity against subtitle space and audience comprehension.

### 3.2 Pedersen's Typology of ECR Translation Strategies

AVT has provision for extralinguistic cultural references (ECRs) using a model drawn up by Pedersen. His typology includes retention, specification, direct translation, generalization, substitution, and omission. Selection of strategy by the subtitler depends on a number of factors including the cultural knowledge of the audience,

technical constraint, and the judgement of the subtitler. Corpus tools intersecr with this typology in offering empirical evidence of how comparable ECRs have been translated in other subtitles.

### 3.3 Descriptive Translation Studies and Translation Norms

Descriptive Translation Studies (DTS) identify regularities and norms in translated texts. In DTS, translation is viewed as a sociocultural activity influenced by institutional and historical factors. Through corpus analysis, corpus-aided AVT reveals patterns and trends in subtitle translation, supporting DTS. It enables researchers and practitioners to investigate whether certain Beijing-specific terms are routinely omitted, domesticated, or foreignized, and what norms these tendencies reflect.

### 3.4 Integrating Corpus Tools with Translation Theory

Corpus tools bridge theory and practice. They operationalize functionalist decisions, guide ECR strategy selection, and expose prevailing norms, aligning with all three theoretical frameworks discussed. When used critically, they inform decisions and also enhance theoretical reflection, making them indispensable in contemporary AVT that involves culture-loaded terms.

## 4. Applications of Corpus-Aided Translation in AVT

### 4.1 Benefits of Corpus Tools in Subtitle Translation

Additionally, translation corpora provide large-scale language information, which is supported by empirical evidence in this manner. These tools include parallel corpora (source and target texts aligned at the sentence or phrase level), bilingual corpora, and specialized subtitle databases such as OpenSubtitles, the China National Subtitle Corpus (CNSC), and the TED Talks multilingual subtitle archive.

In audiovisual translation (AVT), corpora facilitate informed decision-making by enabling the retrieval of recurring patterns and real-life subtitle strategies. For culture-loaded terms—especially dialectal or region-specific expressions—corpora reveal whether specific strategies (e.g., retention, generalization, or omission) are consistently employed, allowing

translators to benchmark their choices against existing practice. As Brendel and Vela note, this increases consistency, reduces subjectivity, and enhances cultural sensitivity in subtitle workflows [8].

### 4.2 Translating Beijing Cultural Terms through Corpus Evidence

Since subtitle material is infused with Beijing culture, translators will often encounter Hutong, Siheyuan, and LaoBeijing-type terms. These are terms with architectural, historical, and affective connotations that cannot be translated literally. Corpus evidence takes a prominent role in determining contextually suitable renderings.

For instance, Hutong is often retained in Pinyin form within various subtitle corpora as observed from the CNSC and OpenSubtitles corpora. A concordance search of the bilingual subtitle corpus of Kang et al. showed that Hutong emerged over 70% of the time in its Pinyin form with minimal annotation, and this suggests that retention is a widespread approach in subtitle practice [9]. Similarly, Siheyuan is sometimes rendered descriptively as "courtyard house," whereas corpus evidence shows that such an explanation is contextual and most times omitted from fast-paced scenes.

The expression Lao Beijing, mostly used to refer to nostalgia or a traditional lifestyle, is more contentious. Corpus analysis translates it as "old Beijing," "traditional Beijing," or leaves it untranslated. Cultural preservation varies by strategy.

### 4.3 Using Concordance Lines to Identify Translation Patterns

Concordance lines—parallel instances of word usage in context—are at the center of corpus-assisted decision-making. Looking for Siheyuan in subtitle corpora, concordance lines indicate the way the term is used across genres and audiences. In international audience-oriented documentaries, for example, Siheyuan is preceded by on-screen images and utilized as "traditional courtyard residence." Fictional drama often ellipsis or abbreviates the term to "home" due to spatial constraints.

This evidence enables subtitlers to scrutinize the feasibility of various methodologies and adjust them based on audiovisual and audience-oriented requirements. With reference to Zou, concordance-based findings reduce individual judgment burden and improve translation quality

in multimodal environments [10].

## 4.4 Case Studies and Empirical Applications

Various studies have employed corpus-based methods to evaluate the subtitles of culture-specific items. Liu constructed a bilingual subtitle micro-corpus from Chinese political dramas (e.g., *Knockout*) and compared content, dialogue, and sociocultural markers with an emphasis on cultural expression strategies. Corpus application was found to enable the marking of repeated subtitle strategies used for culturally marked terms such as douzhe and naliqu even when there were no direct equivalents.

In another study, Zhang et al. built a tailored parallel corpus from Beijing-set films subtitled by fans [11]. By concordance analysis, they revealed that user-generated subtitle data offered rich but inconsistent translation strategies, demonstrating the promise yet instability of subtitle corpora as a resource.

## 4.5 Summary of Practical Value

Corpus tools provide AVT professionals a solid foundation for dealing with culture-loaded terms in subtitles. From the extraction of comparable examples to the checking of translation norms, the tools inform decisions based on evidence from the real world. However, as the succeeding section illustrates, these tools also have their limitations.

## 5. Limitations of Corpus Use in Translating Culture-Loaded Terms

Nevertheless, corpus-aided translation is prone to limitation in the field of audiovisual translation (AVT), particularly when translating culturally rich content into subtitles. In Beijing culture, limitations in this aspect are the most critical in coping with low-frequency, delicate terms. Timeliness, sparse data, and limited context are three key limitations in translation quality.

### 5.1 Timeliness: Corpus Data Obsolescence

One of the core challenges to corpus-based subtitleling is its timeliness. The majority of subtitle corpora are collected from earlier generations of translated content and may not keep up with evolving language usage, audience preferences, or subtitle norms.

According to Gouleti, rich but inconsistent methods are employed through user-generated subtitles, indicating both the richness and the volatile nature of subtitle corpora [12]. Translators will be likely to undertake flat cultural choices when a corpus is not updated with recent subtitles.

AVT is especially liable to such a problem since platform style and orientations evolve at the same rates. The freshness of subtitle corpora thus needs to be considered when selecting or developing datasets.

### 5.2 The Low-Frequency Problem: Sparse Data for Culture-Specific Terms

Another important corpus limitation is the low frequency problem. Region-specific expressions like douzhe, nianggaoqi, or qiezha usually appear seldom or only in a few genres like period dramas or regional comedies. Therefore, such vocabulary may appear sparingly, and in consequence, not even appear at all in large general corpora.

Low-frequency cultural vocabulary like douzhe will occur no more than once or twice in large-scale bilingual subtitle corpora—and frequently in vague settings—so consistent translation strategies may not be possible to determine. Concordance analysis is ineffective where there is insufficient frequency, and the translator will have to fall back on external references or rely on experiential judgment.

Moreover, machine-learning-based subtitle tools learned from such corpora may be low-resourced for these words because they rarely appear. Linguistic quality and cultural authenticity are equally affected by the absence of representative data.

### 5.3 Contextual Limitations: Loss of Audiovisual and Multimodal Cues

All but the most advanced corpus tools only offer written information without audiovisual and multimodal background, i.e., tone, facial expressions, sounds, and visuals. Lack of cultural context is particularly damaging when the words to be translated have extremely context-dependent meanings.

The term Douzhe, for example, can mean either a fan of Beijing opera, a comic actor, or a generally humorously inclined person, depending on tone and context. In a comedic sketch, it should be taken as a compliment; in a dramatic one, it might have nostalgic or critical undertones. Corpus entries do not usually capture such nuances, so it is easy for translators

to interpret it out of context or oversimplify the meaning.

Douzhe tends to be mistranslated as "funny guy" or "comedian" without a visual context, which fails to serve Beijing opera conventions' cultural specificity.

## 5.4 Implications for Translation Accuracy and Cultural Fidelity

These limitations affect accuracy of translation and cultural fidelity. Outdated data generate anachronistic translation; sparse entries cause overgeneralization or deletion; and the lack of audiovisual information results in misrepresentation of cultural nuances. While corpus tools provide a foundation for consistent and evidence-driven translation, they should be supplemented with other means—e.g., ethnographic knowledge, translator judgment, and audiovisual inspection—to achieve authentically faithful and functionally appropriate subtitles.

Therefore, the use of corpus tools alone is not sufficient to provide a correct representation of Beijing-specific phrases in AVT. This is where future work should be improved: developing genre-specific, modern, and multimodal-aided corpora for AVT training.

## 6. Discussion and Implications

The above discussion has revealed that corpora possess significant strengths as well as some flaws when dealing with culture-loaded words, particularly those involving Beijing-local settings. On the one hand, corpora facilitate data-driven, pattern-driven decision-making and provide standards of published translation norms. By providing examples of how terms like Hutong, Siheyuan, and Lao Beijing have previously been translated, they support translator consistency and increase efficiency. High subtitle throughput is useful for multilingual TV productions as well as video streaming services.

However, as illustrated in the above sections, corpus tools are also lacking in several essential ways. They are often lacking in timeliness, especially where corpora are updated to respond to evolving AVT approaches or changing linguistic conventions. Corpus tools are also incapable of capturing non-verbal and multimodal features essential to adequate subtitle translation, including tone of voice, visual cues, and cultural references essential to

performance.

These benefits and drawbacks have an impact on translator training. Corpus literacy should be a core component of AVT training, but it should also be with the addition of critical reflection skills and cross-cultural awareness. Trainees need to be taught to view corpus results as not absolute answers but as one amongst many resources. Pedagogical models need to emphasize comparative analysis of concordance lines, fusion of audio-visual contexts, and awareness for sociolinguistic variation in regional dialects.

Corpus tools may be employed to enhance subtitle production workflows. However, to overcome their shortcomings, they must be combined with ethnographic research, including native speaker interviews and cultural consultant feedback, as well as target audience focus group reviews for the evaluation of interpretability and cultural salience. These platforms can coordinate visual, audio, and textual data to cue translators about multimodal signals better. In general, corpus-aided translation can only be regarded as foundation and not solution.

## 7. Conclusion

The current research explored the role of corpus-assisted translation in Beijing cultural expressions in subtitles, particularly the culture-loaded terms. It highlights the importance of corpus tools in providing empirical support, identifying translation norms, and guiding unchanging and efficient decision-making in audiovisual translation (AVT).

However, the findings also indicate the drawback of corpus tools when applied to culture-specific and infrequent words. It is difficult to translate context-sensitive, nuanced words like douzhe via corpus-based approaches due to outdated corpora, scarce data, and absence of audiovisual context. Such drawbacks can compromise a subtitle's cultural authenticity and functionality.

Given these constraints, human judgment is still necessary. Translators must combine cultural insight, critical consciousness, and contextual awareness into the task of translation—faculties corpus tools cannot replicate. The research encourages the development of new, dynamic, and specialized subtitle corpora that reflect recent trends in translation and incorporate multimodal signals.

Future progress has to integrate corpus technologies with ethnographic responses, focus group testing, and AI-enabled tools that can handle audiovisual environments. It is only through this integration that AVT can continue to maintain cultural richness, linguistic precision, and global availability amidst increasingly advanced translation demands.

## References

[1] Laviosa S, Falco G. Corpora and translator education: Past, present, and future. //Pan J, Laviosa S, eds. Corpora and Translation Education. Singapore: Springer, 2023. (New Frontiers in Translation Studies).

[2] Li D, Corbett J, eds. The Routledge Handbook of Corpus Translation Studies. London: Routledge, 2024. DOI: 10.4324/9781003184454.

[3] Pavesi M. Corpus-based audiovisual translation studies: Ample room for development. //Pérez-González L, ed. The Routledge Handbook of Audiovisual Translation. London & New York: Routledge, 2019: 315–332.

[4] Pedersen J. Subtitling norms for television: An exploration focusing on extralinguistic cultural references. Amsterdam: John Benjamins, 2011. (Benjamins Translation Library, Vol. 98). DOI:10.1075/btl.98.

[5] Yu D. Comparative study of Beijing dialect translation in Teahouse from the aesthetic perspective. English Language Teaching and Linguistics Studies, 2024, 6(4). DOI:10.22158/eltls.v6n4p260.

[6] Yao B, Jiang M, Bobinac T, et al. Benchmarking machine translation with cultural awareness. arXiv preprint, 2023. Available: https://arxiv.org/abs/2305.14328.

[7] Vermeer H J, Reiss K. Towards a General Theory of Translational Action: Skopos Theory Explained. Translated by Christiane Nord. London: Routledge, 2013. DOI: 10.4324 /9781315759715.

[8] Brendel J, Vela M. Quality assessment of subtitles – Challenges and strategies. In: Text, Speech and Dialogue: 25th International Conference, TSD 2022, Proceedings. Lecture Notes in Computer Science, Vol. 13502. Springer, 2022: 52–63. DOI: 10.1007 /978-3-031-16270-1_5.

[9] Kang L, Huang L, Peng N, et al. BigVideo: A large-scale video subtitle translation dataset for multimodal machine translation. //Findings of the Association for Computational Linguistics: ACL 2023. 2023: 5201–5216. Available: https://arxiv.org/abs/2305.18326.

[10] Zou D. Multi-dimensional consideration of cognitive effort in translation and interpreting process studies. //Campbell J, et al., eds. Proceedings of the 15th Biennial Conference of the Association for Machine Translation in the Americas: Users & Providers Track. AMTA, 2022: 416–426. Available: https://aclanthology.org/2022.amta-upg.29/.

[11] Zhang Y, Mahfoodh O H A, Tan D A L. A Corpus-based Evaluation of English-Mandarin Cultural References between Fansubbing and Official Subtitling. Journal of Intercultural Communication, 2024, 24(1): 109–119. DOI:10.36923/jicc.v24i1.291.

[12] Gouleti K. Bilingual subtitling in streaming media: Pedagogical implications. Target: International Journal of Translation Studies, 2023, 35(3): 354–377. DOI:10.1075/target.00013.gou.