# Decoding the Bronze Age: A Comprehensive Technical Report on the Deep Past Initiative Machine Translation Challenge

## 1. Executive Summary

The "Deep Past Challenge: Translate Akkadian to English" represents a frontier task in Natural Language Processing (NLP), situated at the intersection of low-resource Neural Machine Translation (NMT), computational philology, and historical data reconstruction. Hosted by the Deep Past Initiative on Kaggle, this competition challenges participants to translate transliterated Old Assyrian texts—primarily business letters, contracts, and administrative records from the Kanesh trade colony (c. 19th century BCE)—into modern English.[1]

While the surface objective is straightforward—mapping source text to target text—the underlying complexity makes this a distinct "hard problem" in NLP. Unlike high-resource translation tasks (e.g., WMT English-German) where aligned sentence pairs number in the millions, this competition provides a highly constrained, noisy, and structurally mismatched dataset. The training data is provided at the *document level* (full clay tablets), while the test set and evaluation metric operate at the *sentence level*. This discrepancy creates a fundamental data engineering bottleneck that, if unaddressed, caps model performance regardless of architectural scale.[3]

Our comprehensive analysis indicates that the winning strategy is not merely architectural but primarily philological and structural. The dominant architectural paradigm is **ByT5 (Byte-Level T5)**, which bypasses the tokenization failures inherent in standard Transformers when applied to the polyvalent and morphologically rich Akkadian script.[5] However, the decisive competitive advantage lies in the rigorous exploitation of auxiliary metadata—specifically the Sentences_Oare_FirstWord_LinNum.csv file—to reconstruct a precise sentence-level training corpus from the raw document dumps.

This report provides an exhaustive technical roadmap for the competition. It dissects the mathematical properties of the **Geometric Mean (BLEU $\times$ chrF++)** metric, proving why it necessitates a "conservative" translation style over a "fluent" one. We present a robust, mathematically grounded preprocessing pipeline, a hyper-optimized training regimen for ByT5, and an ensemble inference strategy designed to maximize the micro-averaged metrics used for the final leaderboard ranking.

# 2. Domain Analysis: Philology Meets Computation

To engineer an effective translation system, one must first understand the specific linguistic and material properties of the source data. The "Deep Past" here refers to the Old Assyrian period, specifically the archives of Kanesh (modern Kültepe, Turkey).

## 2.1 The Linguistic Landscape of Old Assyrian

Akkadian is an extinct East Semitic language. The dialect in question, Old Assyrian, is distinct from the later Neo-Assyrian or Babylonian dialects often found in broader corpora.[7]

**Morphological Complexity:**

Akkadian utilizes a root-and-pattern morphology (like Arabic or Hebrew). A root consisting of three consonants (radicals) is modified by vowel patterns and affixes to generate meaning.

- *Example:* The root *Š-Q-L* (to weigh/pay) can become *išqul* (he weighed), *išaqqal* (he will weigh), or *šuqulum* (to weigh).
- *Computational Implication:* Standard subword tokenizers (WordPiece, SentencePiece) often fragment these roots destructively. A tokenizer might split *išaqqal* into i + ša + qq + al based on frequency, obscuring the root *Š-Q-L*. This necessitates models that can "see" the sub-token morphological structure, favoring character-level or byte-level approaches.[5]

**The Writing System: Cuneiform to Transliteration:**

The models do not process cuneiform wedges ( $\unicode{x12000}$ ). They process **transliteration**, a scholarly Latinization of the signs. This introduces specific "noise" features:

1. **Polyvalence:** A cuneiform sign can represent a syllable (e.g., *an*), a logogram (a whole word, usually Sumerian, e.g., *DUMU* for "son"), or a determinative (a silent classifier, e.g., *{d}* for god).
2. **Logograms:** These are often written in capital letters in transliteration (e.g., KÙ.BABBAR for silver). The model must learn that KÙ.BABBAR translates to "silver" or "money," not a proper noun.
3. **Determinatives:** Marked with braces {} or superscripts. These are semantic cues not pronounced in speech. For instance, {d}šamaš helps the reader know *Šamaš* is a god. An NMT model must decide whether to treat these as noise or valuable semantic features.[8]

## 2.2 The Nature of the Corpus: "Ancient Excel Sheets"

The dataset consists largely of business correspondence and administrative lists.

- **Formulaic Syntax:** "Thus says A to B: regarding the silver..."
- **Inventory Lists:** "X minas of silver, Y shekels of gold."
- **Implication for Modeling:** The high repetition of formulaic phrases suggests that

models with strong memorization capacities (like Transformers) will perform well, provided they are not hallucinating. However, the *Metric Trap* (discussed in Section 4) punishes synonyms. If the reference translation is "weigh out" and the model predicts "pay," the BLEU score suffers despite semantic correctness.[9]

# 3. Data Ecosystem and Engineering Constraints

The competition provides a fragmented data landscape. Success depends entirely on the ability to synthesize these fragments into a coherent training set.

## 3.1 Dataset Inventory and Structure

| File Name | Granularity | Count | Content Description | Utility |
|---|---|---|---|---|
| train.csv | Document | ~1,561 | Columns: oare_id, transliteration, translation. Contains full texts of tablets. | **Primary Source**, but requires slicing. |
| test.csv | Sentence | ~4,000 | Columns: id, text_id, transliteration. Target for submission. | **Inference Target**. Defines the domain. |
| Sentences_Oare...csv | Metadata | ~9,782 | Columns: text_uuid, first_word_number, translation. Maps sentences to locations in docs. | **The Rosetta Stone** for alignment. |
| publications.csv | Document | ~900 | OCR output from PDFs. | **High Noise**. Use with |

| | | | | extreme caution. |
|---|---|---|---|---|
| eBL_Dictionary.csv | Lexicon | N/A | Dictionary of Akkadian words. | **Auxiliary**. Good for post-processing/lookup. |

## 3.2 The Alignment Bottleneck

The critical pitfall is the mismatch between train.csv (documents) and test.csv (sentences). Training on full documents and predicting sentences yields poor results because the model learns to generate massive blocks of text rather than concise translations.

**The "OARE" Solution:** The file Sentences_Oare_FirstWord_LinNum.csv is the key to unwinding this complexity. It provides the first_word_number for each sentence within a text.[4]

- *Mechanism:* By linking train.csv and the OARE file via oare_id/text_uuid, one can algorithmically "slice" the long transliteration strings in train.csv into individual sentences.
- *Impact:* This transforms the dataset from ~1,500 documents to ~9,700 aligned sentence pairs, a massive increase in supervision signal.[9]

## 3.3 Data Quality and "Noise"

The source text is riddled with scholarly annotations that act as noise for NMT:

1. **Gaps:** Missing signs are represented variously as [x], ..., x x, or (x). These indicate physical damage to the tablet.
2. **Editorial Marks:** Signs like !, ?, /, and * are modern editorial judgments (e.g., "this sign is damaged but readable as...").
3. **Normalization Discrepancy:** The test set uses h and H for the velar fricative, while training data often uses ḫ and Ḫ. Failure to normalize this results in a severe penalty in the chrF++ metric, which is character-sensitive.[3]

**Strategic Insight:** "Clean" training data is a myth in this domain. The winning approach involves *standardizing* the noise (e.g., mapping all gap markers to a single <gap> token) rather than removing it, as the model needs to learn to predict "I cannot read this" when it encounters a gap, rather than hallucinating text.[11]

# 4. The Mathematics of Evaluation

The competition uses a composite metric: the **Geometric Mean of BLEU and chrF++**,

computed via **micro-averaging**.[1]

## 4.1 The Formula and its Derivatives

$$Score = \sqrt{\text{BLEU}_{\text{micro}} \times \text{chrF} + +_{\text{micro}}}$$

**Why Geometric Mean?** The geometric mean ( $\sqrt{A \times B}$ ) is more sensitive to low values than the arithmetic mean. If either BLEU or chrF++ is near zero, the total score collapses. This forces models to be competent in *both* word-level precision (BLEU) and character-level morphology (chrF++). You cannot "game" the system by optimizing only one.[13]

**The Micro-Average Factor:**

- *Macro-average:* Calculate score for each sentence, then average the scores.
- *Micro-average:* Sum the numerators (matches) and denominators (total lengths) across the *entire corpus* first, then calculate the ratio.
- *Implication:* Longer sentences effectively weigh more. A mistake in a long inventory list hurts the score more than a mistake in a short greeting. The "corpus-level" calculation means optimization must focus on maximizing global overlap statistics.[12]

## 4.2 Detailed Metric Analysis

**BLEU (Bilingual Evaluation Understudy):**

BLEU measures the precision of n-grams (1-gram to 4-gram) in the hypothesis against the reference.

- *Strengths:* Checks for fluency and exact vocabulary usage.
- *Weaknesses:* Brittle. It fails on synonyms. If the reference is "silver" and prediction is "money," BLEU gives 0 for that unigram.
- *Relevance:* High BLEU requires the model to memorize the specific "translationese" of the training data (e.g., specific choices of prepositions or verbs used by the OARE annotators).[14]

**chrF++ (Character n-gram F-score):**

chrF++ measures the F-score (harmonic mean of precision and recall) of character n-grams (usually up to 6-grams) and word n-grams (up to 2-grams).

- *Strengths:* Robust to morphology. If the reference is "weighed" and prediction is "weighing", BLEU is 0, but chrF++ gives partial credit for the shared "weigh" character sequence.
- *Relevance:* Critical for agglutinative languages and proper nouns. Akkadian names (e.g., *Aššur-bāni-apli*) have many spelling variations. chrF++ ensures partial credit for getting the name mostly right.[14]

## 4.3 The "Vibes" vs. "Translation" Trap

Analysis of discussion threads highlights a critical strategic pivot.

- **The Trap:** A "good" translation model (e.g., GPT-4) might translate the meaning accurately but use modern phrasing. This results in a low score (~11.8) because it lacks n-gram overlap with the specific, often archaic, reference translations.
- **The Winning Path:** A "parrot" model that memorizes the *patterns* of the training data (e.g., specific ways of listing items) scores higher (~34.2). The metric rewards "pattern matching" over "semantic understanding".[9]
- **Conclusion:** Do not fine-tune on modern English corpora. Fine-tune *only* on the provided domain data to capture the specific "OARE dialect" of English.

# 5. Architectural Imperatives: The Case for Byte-Level Modeling

The leaderboard analysis confirms a stark reality: standard Tokenizer-based models (BERT, mBART, standard T5) are suboptimal for this task. The state-of-the-art (SOTA) solutions are dominated by **ByT5**.[9]

## 5.1 The Tokenization Failure Mode

Standard Transformers use subword tokenizers (WordPiece, SentencePiece) trained on massive corpora (Common Crawl).

- **The Unknown Problem:** When these tokenizers encounter Akkadian transliteration (e.g., *ḫa-muš-tum*), they often fracture it into meaningless fragments or, worse, map unique characters (like ḫ or subscript numbers $_2$) to <unk> (unknown) tokens if those characters weren't in the pre-training corpus.
- **Morphological Blindness:** Splitting a Semitic root like *ptr* (to loosen) based on frequency might separate *p* from *tr*, destroying the semantic unit.

## 5.2 ByT5: The "Token-Free" Solution

**ByT5 (Byte-Level T5)** operates directly on UTF-8 bytes.

- **Mechanism:** It reads the text byte-by-byte. The character "š" is processed as its specific byte sequence (e.g., 0xC5 0xA1).
- **Advantages for Akkadian:**
  1. **No OOV (Out of Vocabulary):** It can process *any* string, no matter how rare or noisy.
  2. **Robustness to Noise:** It treats a typo or an OCR error as a small byte-level perturbation, not a completely different token.[5]
  3. **Granularity:** It allows the model to learn morphology at the finest possible level, essential for the "root-and-pattern" structure of Akkadian.[6]
  4. **Transfer Learning:** ByT5 is pre-trained on the colossal mC4 corpus, giving it a

strong prior on language structure even without specific Akkadian pre-training.

## 5.3 Comparative Analysis: ByT5 vs. Others

Research indicates that while mBART-50 or mT5 perform decently on high-resource languages, ByT5 consistently outperforms them in low-resource (<10k examples) and noisy settings.[16] In the specific context of transliterated ancient languages, ByT5 has shown superior performance in lemmatization and translation tasks compared to standard T5.[17]

**Recommendation:** The architectural backbone of the solution must be google/byt5-base (or small for rapid prototyping). The large variant may be too computationally expensive for the allowed 9-hour runtime without massive gains, given the small dataset size.[1]

# 6. Strategic Implementation Roadmap (Writing Winning Code)

This section details the precise code implementation required to execute the strategy defined above.

## 6.1 Data Loading & Alignment Pipeline

This script is the single most important component. It converts the unusable document data into usable sentence data.

Python

```python
import pandas as pd
import re
import os

# CONFIGURATION
TRAIN_PATH = "/kaggle/input/deep-past-initiative-machine-translation/train.csv"
ALIGNMENT_PATH =
"/kaggle/input/deep-past-initiative-machine-translation/Sentences_Oare_FirstWord_LinNum.csv"

def align_and_slice_data(train_path, alignment_path):
    """
    Core logic to slice document-level transliterations into sentence-level samples.
    """
    print("Loading datasets...")
    df_train = pd.read_csv(train_path)
```

```python
df_align = pd.read_csv(alignment_path)

# Pre-cleaning: Ensure IDs are strings and handle any missing values
df_train['id'] = df_train['id'].astype(str)
# The alignment file links to train via 'text_uuid' or 'oare_id'.
# Inspection of [18] suggests 'oare_id' or 'text_uuid' are keys.
# Let's assume 'text_uuid' maps to 'id' in train based on discussion snippets.

aligned_samples =

# Group alignment metadata by document ID
grouped_align = df_align.groupby('text_uuid')

print(f"Processing {len(grouped_align)} documents for alignment...")

for text_id, group in grouped_align:
    # Find the full document text
    doc_row = df_train[df_train['oare_id'] == text_id] # Adjust key based on exact csv header
    if doc_row.empty:
        # Fallback check for alternate ID columns if necessary
        continue

    full_transliteration = doc_row.iloc['transliteration']

    # Split full text into tokens.
    # CRITICAL: This split must match the logic used to generate 'first_word_number'.
    # Usually simple whitespace splitting is sufficient for transliteration.
    tokens = full_transliteration.split()

    # Sort sentences by their position
    group = group.sort_values('first_word_number')

    for i in range(len(group)):
        row = group.iloc[i]

        # Get start index (0-based conversion if necessary)
        # Snippet  says "max(0, first_word_num - 1)" suggesting 1-based index in CSV.
        start_idx = max(0, int(row['first_word_number']) - 1)

        # Determine end index
        if i < len(group) - 1:
            next_row = group.iloc[i + 1]
            end_idx = max(0, int(next_row['first_word_number']) - 1)
```

```python
        else:
            end_idx = len(tokens)

        # Slice the tokens
        sentence_tokens = tokens[start_idx:end_idx]
        source_text = " ".join(sentence_tokens)
        target_text = row['translation']

        if len(source_text) > 0 and isinstance(target_text, str):
            aligned_samples.append({
                'source': source_text,
                'target': target_text,
                'doc_id': text_id
            })

    print(f"Successfully aligned {len(aligned_samples)} sentence pairs.")
    return pd.DataFrame(aligned_samples)

# Generate the dataset
df_aligned = align_and_slice_data(TRAIN_PATH, ALIGNMENT_PATH)
```

*Rationale:* This code directly implements the "hybrid sentence-level alignment" discussed in the forums.[4] It solves the document-vs-sentence mismatch.

## 6.2 Preprocessing and Normalization

Normalization is critical for the chrF++ score.

Python

```python
def preprocess_text(text, is_source=True):
    if pd.isna(text): return ""
    text = str(text)

    if is_source:
        # 1. Normalize 'h' characters
        # Test set uses 'h/H', train often uses 'ḫ/Ḫ'
        text = text.replace('ḫ', 'h').replace('Ḫ', 'H')

        # 2. Normalize subscripts
```

```python
        # Flatten u₂ -> u2 to match tokenizer/byte consistency
        sub_map = str.maketrans("₀₁₂₃₄₅₆₇₈₉", "0123456789")
        text = text.translate(sub_map)

        # 3. Standardize Gaps [12]
        # Map all variations of missing text to a single token
        text = re.sub(r'\[x\]|\(x\)|\.\.\.|x\s?x', ' <gap> ', text)
        text = re.sub(r'\s+', ' ', text).strip()

    else:
        # Target Normalization (English)
        # Remove editorial marks that might confuse the model
        text = text.replace('?', '').replace('!', '')
        text = re.sub(r'\s+', ' ', text).strip()

    return text

# Apply to DataFrame
df_aligned['source_clean'] = df_aligned['source'].apply(lambda x: preprocess_text(x,
is_source=True))
df_aligned['target_clean'] = df_aligned['target'].apply(lambda x: preprocess_text(x,
is_source=False))
```

## 6.3 Metric Implementation (GeoMean)

This function mimics the official Kaggle evaluation.

Python

```python
import numpy as np
import sacrebleu

def compute_metrics(eval_preds):
    preds, labels = eval_preds

    # Decode generated bytes to text
    decoded_preds = tokenizer.batch_decode(preds, skip_special_tokens=True)
    decoded_labels = tokenizer.batch_decode(labels, skip_special_tokens=True)

    # Post-process for metric calculation (strip whitespace)
```

```python
    decoded_preds = [pred.strip() for pred in decoded_preds]
    decoded_labels = [[label.strip()] for label in decoded_labels] # Nested list for sacrebleu

    # BLEU (Micro-average)
    bleu = sacrebleu.corpus_bleu(decoded_preds, decoded_labels)

    # chrF++ (Micro-average)
    chrf = sacrebleu.corpus_chrf(decoded_preds, decoded_labels, word_order=2)

    # Geometric Mean
    # Note: sacrebleu returns score out of 100.
    geo_mean = np.sqrt(bleu.score * chrf.score)

    return {
        "bleu": bleu.score,
        "chrf": chrf.score,
        "geo_mean": geo_mean
    }
```

*Note on word_order=2:* This parameter in corpus_chrf enables the "++" part of chrF++, which includes word bigrams in the calculation, crucial for correlation with human judgment.[14]

## 6.4 Model Training Loop (Hugging Face Trainer)

We use the Seq2SeqTrainer for optimized training.

Python

```python
from transformers import AutoTokenizer, AutoModelForSeq2SeqLM,
Seq2SeqTrainingArguments, Seq2SeqTrainer

MODEL_CHECKPOINT = "google/byt5-base"

# Load Tokenizer & Model
tokenizer = AutoTokenizer.from_pretrained(MODEL_CHECKPOINT)
model = AutoModelForSeq2SeqLM.from_pretrained(MODEL_CHECKPOINT)

# Tokenization Function
def tokenize_function(examples):
    model_inputs = tokenizer(examples["source_clean"], max_length=512, truncation=True)
```

```python
    with tokenizer.as_target_tokenizer():
        labels = tokenizer(examples["target_clean"], max_length=512, truncation=True)
    model_inputs["labels"] = labels["input_ids"]
    return model_inputs

# Convert to HF Dataset
from datasets import Dataset
hf_dataset = Dataset.from_pandas(df_aligned[['source_clean', 'target_clean']])
# Split (Stratified recommended, random shown for brevity)
hf_dataset = hf_dataset.train_test_split(test_size=0.1)
tokenized_datasets = hf_dataset.map(tokenize_function, batched=True)

# Training Arguments
args = Seq2SeqTrainingArguments(
    output_dir="byt5-akkadian-v1",
    evaluation_strategy="epoch",
    learning_rate=2e-4,          # Higher LR for ByT5
    per_device_train_batch_size=8,
    per_device_eval_batch_size=8,
    weight_decay=0.01,
    save_total_limit=3,          # Save space
    num_train_epochs=15,         # Long training needed
    predict_with_generate=True,
    fp16=True,                   # Mixed precision for speed
    generation_max_length=512,
    metric_for_best_model="geo_mean",
    greater_is_better=True,
)

trainer = Seq2SeqTrainer(
    model=model,
    args=args,
    train_dataset=tokenized_datasets["train"],
    eval_dataset=tokenized_datasets["test"],
    tokenizer=tokenizer,
    compute_metrics=compute_metrics,
)

# trainer.train()
```

# 7. Hyperparameter Optimization and Training

# Dynamics

The difference between a baseline model (34.0) and a competitive model (36.0+) lies in the hyperparameter tuning.

## 7.1 The Search Space

Based on community findings and ByT5 literature, the following grid yields optimal results.[5]

| Hyperparameter | Value | Justification |
| --- | --- | --- |
| Learning Rate | 1e-4 - 3e-4 | ByT5 is pre-trained with a high LR; fine-tuning requires maintaining this momentum to adjust byte embeddings. |
| Batch Size | 8 | constrained by GPU memory (T4 on Kaggle). ByT5 sequences are 4x longer than subword sequences. |
| Max Source Length | **512** | Progression from 385 -> 475 -> 512 showed steady gains. Akkadian sentences can be verbose. |
| Label Smoothing | 0.1 | Prevents the model from becoming overconfident, which aids generalization on the "brittle" BLEU metric. |
| Epochs | 10 - 20 | Byte-level models converge slower than token-level models. |
| Beam Size | 8 | High beam width allows the model to explore more potential translation paths, |

| | | critical for finding the exact phrasing required by BLEU. |
|---|---|---|
| **Length Penalty** | 1.096 | Slightly encourages longer output. The reference translations often include explanatory phrases. |

## 7.2 Regularization and Overfitting

With only ~9,700 sentences, overfitting is a major risk.

- **Weight Decay:** Set to 0.01 to penalize large weights.
- **Early Stopping:** Monitor eval_geo_mean. Stop if no improvement after 3 epochs.
- **Dropout:** Standard Transformer dropout (0.1) is usually sufficient.

# 8. Inference, Post-Processing, and Ensembling

Inference is not just calling model.generate(). It is a multi-step pipeline designed to maximize the metric.

## 8.1 Chunked Beam Search

Given the 9-hour runtime limit on Kaggle, processing the 4,000 test sentences with a beam size of 8 can be slow.

- **Strategy:** Implement "Chunked Beam Search".[13] Process the test set in batches (chunks) to manage memory and allow for intermediate saving. If the kernel crashes, you don't lose all progress.

## 8.2 The Power of Ensembling

The top scores on the leaderboard (~36.5) are undoubtedly ensembles.

- **Method:** Train 3-5 ByT5 models.
  - *Model A:* Trained on 100% train.csv.
  - *Model B:* Trained on train.csv + ORACC external data.
  - *Model C:* Trained with different random seed.
- **Inference:**
  - Generate predictions from all models.
  - **Weighted Voting:** If the models are diverse, use "Centroid Selection" (select the hypothesis that minimizes distance to others) or simply average the probabilities at each decoding step (Ensemble Decoding). Given the complexity of implementing ensemble decoding in Hugging Face, **System Combination** (e.g., using a simple

voting mechanism or choosing the output with the highest average length/confidence) is a practical alternative.

## 8.3 Post-Processing Rules

Before submission, run the predictions through a final cleaning filter:

1. **Collapse Gaps:** re.sub(r'(<gap>\s*)+', '<gap> ', text) - prevent the model from stuttering "………".
2. **Fix Punctuation:** Ensure exactly one space after commas and periods. The metric tokenization relies on standard spacing.
3. **Capitalization Check:** Ensure the first letter is capitalized (English convention), as references likely are.

# 9. Synthesis and Strategic Recommendations

The **Deep Past Challenge** is a masterclass in "Data-Centric AI." The models (ByT5) are standard, but the data engineering is bespoke.

**Key Takeaways:**

1. **Philology First:** Understanding the Sentences_Oare file is more valuable than tweaking the Transformer architecture. The document-to-sentence alignment is the primary lever for score improvement.
2. **Byte-Level Supremacy:** The linguistic properties of Akkadian (morphology, transliteration noise) make ByT5 the only viable architectural choice.
3. **Metric Awareness:** The Geometric Mean of Micro-Averaged BLEU/chrF++ demands a translation style that is conservative, consistent, and strictly adheres to the training data's "dialect" of English.
4. **Normalization is Non-Negotiable:** Failure to map ḫ to h or normalize gaps will cap performance at the baseline level.

**Final Recommendation:**

Begin by implementing the alignment script provided in Section 6.1. Train a byt5-small model to verify the pipeline. Once valid, scale to byt5-base, implement the ensemble strategy, and rigorously cross-validate using a stratified split that respects document boundaries. This path leads directly to the top tier of the leaderboard.

---

## Deliverables Checklist

- [x] **Report:** Comprehensive analysis of task, data, and strategy.
- [x] **Code:** Full Python snippets for Alignment, Preprocessing, Metric, and Training.
- [x] **Strategy:** ByT5 Ensembling + OARE Alignment.

- [x] **Metric Analysis:** GeoMean breakdown and "Vibes" warning.
- [x] **Data Analysis:** Detailed breakdown of file utility and alignment logic.

## Works cited

1. Deep Past Challenge - Translate Akkadian to English - Kaggle, accessed February 12, 2026, https://www.kaggle.com/competitions/deep-past-initiative-machine-translation
2. Andrew A.N. Deloucas (@aandeloucas.com) — Bluesky, accessed February 12, 2026, https://bsky.app/profile/aandeloucas.com
3. Deep Past Challenge - Translate Akkadian to English - Kaggle, accessed February 12, 2026, https://www.kaggle.com/competitions/deep-past-initiative-machine-translation/discussion/667065
4. Deep Past Challenge - Translate Akkadian to English - Kaggle, accessed February 12, 2026, https://www.kaggle.com/competitions/deep-past-initiative-machine-translation/discussion/671781
5. Proceedings of the 1st Workshop on Machine Learning for Ancient Languages (ML4AL 2024) - ACL Anthology, accessed February 12, 2026, https://aclanthology.org/2024.ml4al-1.pdf
6. Are Character-level Translations Worth the Wait? Comparing ByT5 and mT5 for Machine Translation | Transactions of the Association for Computational Linguistics - MIT Press Direct, accessed February 12, 2026, https://direct.mit.edu/tacl/article/doi/10.1162/tacl_a_00651/120650/Are-Character-level-Translations-Worth-the-Wait
7. Deep Past Challenge - Translate Akkadian to English - Kaggle, accessed February 12, 2026, https://www.kaggle.com/competitions/deep-past-initiative-machine-translation/discussion/663357
8. From Clay Tablets to Kaggle. Why This Akkadian Dataset Is Hard ..., accessed February 12, 2026, https://mhr007.medium.com/from-clay-tablets-to-kaggle-1f8905cd1804
9. Deep Past Challenge - Translate Akkadian to English | Kaggle, accessed February 12, 2026, https://www.kaggle.com/competitions/deep-past-initiative-machine-translation/discussion/672511
10. Deep Past Challenge - Translate Akkadian to English - Kaggle, accessed February 12, 2026, https://www.kaggle.com/competitions/deep-past-initiative-machine-translation/data
11. BYT-Ensemble | Script - Kaggle, accessed February 12, 2026, https://www.kaggle.com/code/anthonytherrien/byt-ensemble-script
12. Akkadian MT: Preprocessing, Ensemble Inference - Kaggle, accessed February 12, 2026,

https://www.kaggle.com/code/carp0308/akkadian-mt-preprocessing-ensemble-inference

13. ByT5 + Optuna Tuning + Chunked Beam Search - Kaggle, accessed February 12, 2026, https://www.kaggle.com/code/harukiharada/byt5-optuna-tuning-chunked-beam-search

14. The AI community building the future. - Hugging Face, accessed February 12, 2026, https://huggingface.co/metrics

15. Deep Past Challenge - Translate Akkadian to English | Kaggle, accessed February 12, 2026, https://www.kaggle.com/competitions/deep-past-initiative-machine-translation/discussion/668327

16. A Few Thousand Translations Go A Long Way! Leveraging Pre-trained Models for African News Translation - Microsoft, accessed February 12, 2026, https://www.microsoft.com/en-us/research/wp-content/uploads/2022/05/a_few_thousand_translations_go.pdf

17. Other Workshops and Events (2025) - ACL Anthology, accessed February 12, 2026, https://aclanthology.org/events/ws-2025/