

Model Selection Report for PDF Translator – Language Translation Module

Objective:

To identify and select an optimal **language translation model** for a web-based **Hindi ↔ English PDF Translator**, considering factors like cost, deployment feasibility, performance, and scalability.

Evaluated Approaches

1. Transformer Models via Hugging Face API

- **Models Tried:**
 - Helsinki-NLP/opus-mt-en-hi
 - Helsinki-NLP/opus-mt-hi-en
 - ai4bharat/indicTrans2
- **Pros:**
 - Reliable translation for Indian languages
 - Easy to integrate
- **Cons:**
 - Hugging Face inference API has strict **rate limits** (for free-tier)
 - IndicTrans2 requires **GPU** or powerful local setup for smooth usage

2. Multilingual Model: mbart-large-50-many-to-many-mmt

- **Pros:**
 - Supports **both Hindi and English** in a single model
 - Handles multilingual scenarios well
 - Easy to deploy via Hugging Face Transformers locally or on lightweight environments
- **Cons:**
 - Slightly heavier than single-direction models
 - Moderate inference speed without GPU

3. Deep Translator Libraries (e.g., Google, LibreTranslate)

- **Pros:**
 - Abstracts underlying translation engine
 - Quick to use for short texts
- **Cons:**
 - **Google Translate** is a **paid service** beyond trial
 - **LibreTranslate** is **limited in language quality**
 - Rate limits or unstable for production use without self-hosting

4. Paid API-based Models (e.g., OpenAI, Google Translate API)

- **Pros:**
 - Highly fluent, natural translations
 - Context-aware and grammar-correct
- **Cons:**
 - **Paid access only**
 - Not viable for free-tier or open-source constrained projects

5. Open-Source LLMs (e.g., LLaMA, Mistral via Ollama or HF Inference)

- **Pros:**

- Highly customizable via prompts
- Can be tuned for intelligent translation rules (e.g., skip acronyms)
- **Cons:**
 - Require **large downloads** (10–20 GB)
 - Need **local GPU** or **paid inference endpoint**
 - Hugging Face inference often limited under free-tier

Final Model Selection:

facebook/mbart-large-50-many-to-many-mmt

- **Reason for selection:**
 - Open-source
 - Multilingual (bidirectional support in a single model)
 - Can run locally on CPU with moderate performance
 - No dependency on paid APIs or large LLMs
 - Compatible with Hugging Face Transformers for deployment in Streamlit-based apps

For High-Accuracy Option (LLMs with Prompt Templating)

- For **maximum translation accuracy** and **fine-grained control**, advanced LLMs like **GPT-4, Claude, or Gemini Pro** can be used with **prompt templating and preprocessing-aware translation**.
- These models can:
 - Understand and **preserve technical terms**
 - Skip **acronyms and capitalized words** intelligently
 - Maintain **natural sentence structure and flow**
- However, these options:
 - Are **fully paid APIs**
 - Require **usage-based billing** or enterprise-level access
 - May need **pre-tokenized input** for long documents

Conclusion:

- After evaluating multiple translation model options across deployment cost, infrastructure requirement, and multilingual capability, we have selected **mBART-50** as the translation engine for this project. It meets the criteria of **free, multilingual, deployable, and reliable**, making it an ideal choice in a resource-constrained environment.
- For future upgrades, integrating paid LLMs like GPT can significantly enhance translation quality through prompt-level preprocessing and intelligent handling.

Preprocessing Role in Translation

When using LLMs:

- Preprocessing logic can be embedded inside the **prompt template**.
- Example:

"Translate the paragraph into Hindi. Do not translate abbreviations like AI, NASA or words in full capital letters like PDF, ML. Preserve formatting and symbols."

-  Future tuning is simple — just update the prompt.
- No need for code-based filtering unless using external regex for edge cases.

When not using LLMs:

- Need to build a **manual preprocessing layer in code**.
- Ensures translation models like mBART or Helsinki avoid:
 - Mis-translating acronyms
 - Breaking formatting
 - Altering symbols and URLs
- Requires both **preprocessing (before translation)** and **postprocessing (after translation)** stages.

Preprocessing Tasks & Rules

1. Preserve Abbreviations & Acronyms

Use a predefined **abbreviation dictionary**:

ABBREVIATIONS = {

```
'AI', 'ML', 'API', 'URL', 'PDF', 'HTML', 'CSS', 'JS', 'SQL', 'JSON',
'HTTP', 'HTTPS', 'NASA', 'FBI', 'CEO', 'CTO', 'PhD', 'MBA', 'USA',
'UK', 'UAE', 'CPU', 'GPU', 'RAM'
```

}

- Replace each abbreviation with a placeholder before translation.
- Restore after translation.

2. Block Patterned Phrases

Prevent translation of specific patterns using regex:

- URLs: http://, https://, www.
- Emails: @, .com
- File names: .pdf, .png, FILE_NAME_123.PDF
- Versions: 1.0.2a, v2.3.1-beta

All such tokens will be masked before translation and restored afterward.

3. Preserve FULL CAPS Words

Identify and skip all-uppercase words via regex:

- Replace them with placeholders like __CAPS_PDF__
- Restore post-translation.

4. Modernize Hindi Terminology

Use a **Hindi normalization dictionary** to modernize outdated or formal Hindi words.

{

```
"नमूना": "सैंपल",
```

"सूचना": "इन्फोर्मेशन",
"संगणक": "कंप्यूटर",
"अंतर्जाल": "इंटरनेट"

}

- Applied as a final post-processing step when **Hindi is the output language**

5. Preserve Numbers, Symbols, and Punctuation

Do not alter:

- Decimal numbers (e.g., 12.5)
- Punctuation (.,?;:')
- Currency/units (₹, \$, %)
- Technical symbols (=, +, <, >, #)

Ensure these are passed through untouched or wrapped in tokens during translation.

6. Lowercase Handling

- Input text may be lowercased (optional) if the model performs better on normalized text.
- However, abbreviations and formatting tokens must remain intact.

📌 Summary

Preprocessing Rule	Applied When LLM	Applied When No LLM
Prompt-based filtering	<input checked="" type="checkbox"/> Yes	<input type="checkbox"/> No
Code-based masking (acronyms)	<input type="checkbox"/> No	<input checked="" type="checkbox"/> Yes
Regex filtering (URLs, versions)	<input type="checkbox"/> No	<input checked="" type="checkbox"/> Yes
Modern Hindi substitutions	<input type="checkbox"/> Optional	<input checked="" type="checkbox"/> Yes
Placeholder restoration	<input type="checkbox"/> No	<input checked="" type="checkbox"/> Yes

☒ Final Note:

Preprocessing is **optional but powerful** when using LLMs. However, when relying on **non-LLM transformers like mBART or Helsinki**, it becomes **essential** to ensure intelligent and lossless translations, especially for technical documents.

PDF Formatting & Layout Retention Report – PDF Translator

Objective:

To ensure that the **translated PDF file** retains the **visual structure** and **formatting** of the original input PDF. This includes page layout, fonts, bounding boxes, images, tables, and structural flow.

Formatting Strategy

Your solution is built around **layout-aware text block extraction and reinsertion**, powered by PyMuPDF (fitz) for both reading and writing.

Extracting Format-Aware Blocks (from pdf_reader.py)

- `extract_text_blocks()` retrieves:
 - **Text content**
 - **Bounding box (x0, y0, x1, y1)** for each block
 - **Page number**
 - **Block type** (optional: text/image/table)
- **Images are separately extracted** using `extract_images_info()`
- Output: A list of structured dictionaries that help track text + layout per page

Rebuilding Translated PDF (from pdf_writer.py)

- Translated text is **reinserted at original locations** using:
 - `page.draw_rect(...)` → clears the area
 - `page.insert_htmlbox(...)` → adds styled text preserving font and alignment
- Page dimensions, margins, and fonts are retained
- Default font: "helv" (Helvetica) with fallback to helv-bold or helv-oblique based on flags

Additional Techniques & Enhancements

Technique	Purpose
Sentence Segmentation	Break long paragraphs into shorter, more manageable units for translation
Chunk-wise Translation	Avoids API timeout, memory errors, and supports streaming for large documents
Translation Caching	Skips already translated text blocks to reduce API calls and cost
Context-Aware Glossary	Predefined mappings to consistently translate technical or domain-specific words
Language Detection (optional)	Auto-detect source language if user uploads ambiguous or bilingual files

Specific Implementation Points

1.  **Position Accuracy**
 - Translated text is placed **exactly within original bounding box**
 - Minor height expansion ($y1 \pm= 10$) ensures content fits
2.  **Page Size Preservation**
 - Page width/height extracted from original file

- Used in `create_simple_translated_pdf()` when building new files

3. 🖼 Image Retention

- Image info is extracted, but full image reinsertion logic can be added later (placeholder supported)

4. 🎨 Font & Color Matching

- Font flags are mapped to "helv" variants
- Text color is restored using RGB values derived from original color integer

✓ Summary

Your current system uses a **layout-aware intermediate layer** to extract and rebuild PDF content. Combined with block-level bounding boxes and HTML-based reinsertion, it achieves:

- Preserved **layout and position**
- Accurate **text flow**
- Compatibility with **multi-page** documents

The additional techniques outlined (sentence segmentation, chunking, caching, glossary, detection) further enhance reliability and performance, especially when integrated with external translation APIs.

Streamlit UI Report – PDF Translator

Objective:

To design a **user-friendly and efficient interface** that enables users to:

- Upload a PDF
 - Select the translation direction (Hindi ⇌ English)
 - Trigger intelligent translation
 - Download a translated PDF with preserved formatting
-

Streamlit UI Structure

1. App Configuration

- **Page Title:** PDF Translator
 - **Page Icon:** 
 - **Layout:** wide for spacious UI
-

2. Sidebar – Live Translator Tester

python

CopyEdit

st.sidebar

- **Header:**  *Translation Tester*
 - Allows users to:
 - Choose translation direction (Hindi ⇌ English)
 - Try live filtering with a text input box
 - **Purpose:** Helps test preprocessing and translation logic interactively
-

3. Main UI Flow

◆ Step 1: Upload PDF

- **Component:** st.file_uploader(...)
 - **Validates:**
 - File size (MAX_FILE_SIZE_MB)
 - PDF format (validate_pdf)
 - Text presence (has_extractable_text)
 - **Displays:**
 - PDF metadata (pages, dimensions, size)
 - Status messages (error, warning, success)
-

◆ Step 2: Select Translation Direction

- **Component:** st.selectbox(...)
 - Populates from TRANSLATION_DIRECTIONS dict
 - Shows "From" and "To" languages in two columns for clarity
-

◆ Step 3: Translate PDF

- **Trigger:** st.button("👉 Start Translation")
- **Flow:**
 1. Extract text blocks with coordinates
 2. Clean and flatten text
 3. Translate text blocks via translate_text_blocks(...) with a callback for progress updates
 4. Reconstruct translated PDF:

- Primary method: `create_translated_pdf(...)`
 - Fallback: `create_simple_translated_pdf(...)` if layout fails
 - **Progress Feedback:**
 - Progress bar (`st.progress`)
 - Live status updates (`st.empty()`)
-

◆ Step 4: Download Section

- **Conditional:** Rendered only if `st.session_state["translation_complete"] = True`
 - **Component:** `st.download_button(...)`
 - Allows user to download the translated PDF
 - Option to **restart process** (`st.button("🔄 Translate Another")`) by resetting session state
-

State Management

- Uses `st.session_state` to track:
 - Translation completion
 - Translated PDF binary
 - Uploaded filename
-

User Metrics Displayed

After translation:

- Number of text blocks
 - Total characters translated
 - Final PDF size in KB
-

Summary of Key Features

Section Purpose

Sidebar Tester	Test filtering logic interactively
Upload	Validate PDF and display its metadata
Direction	Let user select translation direction clearly
Translate	Full translation pipeline with progress UI
Download	Enable user to save the translated PDF easily

Final Notes:

Your UI offers a **clean, minimalistic, and functional interface** tailored for non-technical users. It ensures:

- Smooth user experience with visual feedback
- Robust error handling
- Interactive exploration via sidebar
- Modularity for future additions like:
 - Language auto-detection
 - Format preview
 - OCR for scanned/image-based PDFs



Advanced Enhancement – PDF Formatting (Approach 2)

🎯 Objective:

To implement a **layout-preserving, high-fidelity PDF translation system** using deep learning-based **layout analysis, OCR, and advanced translation models**, ensuring that the translated output mirrors the visual and structural design of the original PDF document.

📝 Overall Architecture Workflow

1. Input Processing

- Users upload a PDF via:
 - Streamlit Web UI
-

2. Layout Analysis

- **Model Used:** [DiT \(Document Image Transformer\)](#) by Microsoft (via UniLM)
 - **Purpose:**
 - Analyze and detect the document layout
 - Identify structural elements: text blocks, headers, footers, tables, sections
 - Retain **spatial relationships** between visual components
-

3. OCR (Optical Character Recognition)

- **Engine Used:** [PaddleOCR](#)
 - **Purpose:**
 - Extract **machine-readable text** from rendered PDF images
 - Preserve **bounding box coordinates** for re-insertion
 - **Benefit:**
 - Works well with scanned/image-only PDFs (unlike PyMuPDF)
 - Provides **language-specific OCR pipelines** (including Hindi)
-

4. Translation Engine

- **Default:** Google Translate API (accessible, easy-to-integrate)
 - **Best Quality (Optional):** OpenAI GPT-4 API
 - Configurable via config.yaml
 - **Function:**
 - Receives cleaned text from OCR
 - Applies intelligent translation
 - Optionally supports **prompt-tuned logic** (for GPT) to skip acronyms, preserve formatting
-

5. Layout Reconstruction

- Translated text is:
 - Replaced at the **same bounding boxes** extracted during OCR
 - Ensures formatting, alignment, fonts, and block positioning are retained
 - Uses PDF rendering libraries like PyMuPDF or pdfplumber to rebuild output with original layout fidelity
-

6. Output Generation

- A **translated PDF** is generated:
 - Visually identical in layout to the original

- Text now in target language
 - Output is returned to:
 - Web user (via Streamlit download button)
 - API consumer (as file response)
-

Technical Implementation Workflow

Setup Process

Component	Description
config.yaml	Stores API keys, model paths, translation settings
Model Downloads	Downloads DiT and PaddleOCR models if not cached
GPU Acceleration	Supports CUDA-based PaddleOCR for high-speed inference

Key Advantages

Feature	Benefit
DiT Layout Parsing	Advanced structural understanding of documents
OCR Support	Enables translation of scanned/image-based PDFs
Position-aware Insertion	Maintains document's visual flow
Modular Config	Easy to switch translation backends
GPU-Ready Architecture	High throughput for large files

Future Enhancements (Optional)

- Add support for **table structure detection** using TableNet or Donut
 - Incorporate **font and color preservation** per word
 - Add **real-time preview** of translation in browser (via HTML canvas)
 - Add **text editing overlay** for manual corrections before final PDF export
-

Conclusion

This advanced architecture leverages **modern document layout modeling (DiT)**, **accurate OCR (PaddleOCR)**, and **modular translation engines (Google/GPT)** to deliver **near-original layout fidelity** in translated PDFs. It's ideal for handling:

- Complex layouts
- Scanned files
- High-accuracy enterprise use cases

PDF Translator – Multi-Stage Processing Pipeline (Approach 2)

This document describes the architecture and responsibilities of each processing stage in the pipeline used for intelligent and layout-preserving PDF translation.

Stage 1: PDF Parsing & Intermediate Layer Creation

Technology:

Goal: Extract raw content and preserve spatial structure
Output: Intermediate Layer (IL) JSON

Key Points:

- Parses text, position, fonts, tables, images
- Maintains spatial layout for reassembly

Stage 2: Layout OCR Processing

Technology:

Goal: OCR for scanned/handwritten PDFs

Key Points:

- Adds positional text from OCR layers
- Integrates with existing IL pipeline

Stage 3: Paragraph Recognition

Technology:

Goal: Group text lines into logical paragraphs

Key Points:

- Detects word boundaries using regex (e.g., [0-9A-Za-z]+)
- Avoids broken lines or mid-word splits

Source: Inspired by BabelDOC segmentation strategies

Stage 4: Style and Formula Processing

Technology:

Goal: Preserve formatting and math expressions

Key Points:

- Captures bold, italic, sub/superscript styles
- Preserves LaTeX-style math or in-line formulas
- Generates metadata for reconstruction

Stage 5: Intermediate Layer Translation

Technology:

Goal: Translate content inside the Intermediate Layer

Key Features:

- Works with OpenAI-compatible LLM APIs
- Batch translation support
- **QPS rate limiting** (default: 4 req/sec)
- **Translation caching** to avoid duplication

Stage 6: Typesetting Processing

Technology:

Goal: Reflow translated text while preserving layout

Key Points:

- Handles text expansion/shrinkage due to translation

- Smart positioning and line-breaking
 - Adapts layout dynamically for multi-language content
-

Stage 7: Font Mapping

Technology:

FontMapper

Location:

yadt/document_il/utils/fontmap.py

Goal: Ensure cross-language font compatibility

Key Points:

- Maps original fonts to compatible Hindi/English equivalents
 - Avoids missing glyphs or fallback fonts
 - Enables smooth multilingual rendering
-

Stage 8: PDF Generation

Technology:

PDFCreation

Goal: Generate the final output module

PDF

Output Modes:

- **Dual PDF:** Original + Translation side-by-side or interleaved
- **Mono PDF:** Only translated version
- Supports layout options:
 - Alternating pages
 - Side-by-side columns
 - Bilingual stacking