

1. Chat files -

Data Source	Format	Chunking Method	Rationale	Strategy Name
Chats	.txt	<ul style="list-style-type: none"> Fixed-size chunking with overlap splits where each chunk partially repeats part of the previous one to preserve context. Participant names and their numbers are extracted too. 	Context is not lost at chunk boundaries, and conversational flow is maintained.	Custom (General/Naive chunking strategy of ragflow does not allow chunk overlap)
Chats	.pdf	<ul style="list-style-type: none"> Extract chat text from the PDF and structure it into messages. Split the text into fixed-size chunks with partial overlap to keep context. Identify unique participant names and count them. 	Context is not lost at chunk boundaries, and conversational flow is maintained.	Custom
Chats	.json	<ul style="list-style-type: none"> Dividing each message in json format Grouping 5-10 messages together to form a chunk 	Easy to process for the LLM and chunk size is consistent.	Custom
Chats	.xml	<ul style="list-style-type: none"> Split the chat XML data into chunks by grouping a fixed number of<Message> elements (5-10 messages) together. 	Maintains XML structure, easy for LLM to process and chunk size is consistent.	Custom

2. Bank Statements -

Data Source	Format	Chunking Method	Rationale	Strategy Name
Bank Statements	pdf	<ul style="list-style-type: none">• Use <code>extract_tables</code> from PyMuPDF to extract structured transaction tables from bank statement PDFs.• Group extracted transactions into chunks of 5–10 rows, with some transactions overlapping between adjacent chunks to maintain context.• Attach the column headers to each transaction value in every chunk to ensure clarity and consistency.• Treat non structured data as a separate chunk	Preserves transactional context, and is scalable and robust.	Custom
Bank Statements	xlsx	<ul style="list-style-type: none">• Group transactions into chunks of 5–10 rows each.• Include an overlap of a few transactions• Attach the column headers to each transaction value within every chunk for clarity.	Preserves transactional context and is scalable and robust.	Custom

Bank Statements	xls	<ul style="list-style-type: none"> ● Group transactions into chunks of 5–10 rows each. ● Include an overlap of a few transactions ● Attach the column headers to each transaction value within every chunk for clarity ● Treat non structured data as a separate chunk 	Preserves transactional context and is scalable and robust.	Custom
Bank Statements	csv	<ul style="list-style-type: none"> ● Group transactions into chunks of 5–10 rows each. ● Include an overlap of a few transactions ● Attach the column headers to each transaction value within every chunk for clarity ● Treat non structured data as a separate chunk 	Preserves transactional context and is scalable and robust.	Custom
Bank Statements	html	<ul style="list-style-type: none"> ● Group transactions into chunks of 5–10 rows each. ● Include an overlap of a few transactions ● Attach the column headers to each transaction value within every chunk for clarity ● Treat unstructured data separate chunk 	Preserves transactional context and is scalable and robust.	Custom

3. Video

Data Source	Format	Chunking Method	Rationale	Strategy Name
Video	.video/*	<ul style="list-style-type: none">• First the video can be transcribed using ASR and then• semantic chunking can be done, i.e. similar sentences can be chunked together (sentence transformers)• can add time stamp based overlap	Preserves flow and context, retrieval is good.	Custom

4. Image

Data Source	Format	Chunking Method	Rationale	Strategy Name
Image	.image/*	<ul style="list-style-type: none">• For retrieval, slide fixed size window over image, convert to embeddings and store them• For Ocr, run ocr first and store similar to pdf	Preserves flow and context, retrieval is good.	Custom

5. Audio

Data Source	Format	Chunking Method	Rationale	Strategy Name
Audio	.audio/ *	<ul style="list-style-type: none">• First the video can be transcribed using ASR and then• semantic chunking can be done, i.e. similar sentences can be chunked together (sentence transformers)• can add time stamp based overlap	Preserves flow and context, retrieval is good.	Custom

6. Others

Data Source	Format	Chunking Method	Rationale	Strategy Name
Others	pdf	<ul style="list-style-type: none">• PDF can be parsed into structured data using NER (used for resumes)• Using the lowest section titles as the basic unit for chunking documents. Therefore, figures and tables in the same section will not be separated, which may result in larger chunk sizes. (for manuals)• For research papers, it	<ul style="list-style-type: none">• NER extracts key fields as meaningful chunks for precise data capture.• Lowest section titles group related content, preserving instructional integrity.	<ul style="list-style-type: none">• Resume• Manual• Paper• Laws• Presentation

		<p>can be chunked using sections (eg, abstract, intro)</p> <ul style="list-style-type: none"> • Legal documents are chunked by detecting structured text features like hierarchical headings (e.g., "ARTICLE") as split points, preserving document structure. • In case of presentations, every slide is treated as a separate chunk. 	<ul style="list-style-type: none"> • Results in good info retrieval • Hierarchical headings keep legal provisions intact in chunks. • Each slide stays a clear, self-contained content unit. 	
Others	xlsx	<ul style="list-style-type: none"> • For Q&A documents, each qa pair is treated as one chunk • For table data, each row is considered as a separate chunk 	<ul style="list-style-type: none"> • Context between each pair is preserved 	<ul style="list-style-type: none"> • Q&A • Table
Others	xls	<ul style="list-style-type: none"> • For Q&A documents, each qa pair is treated as one chunk • For table data, each row is considered as a separate chunk 	<ul style="list-style-type: none"> • Context between each pair is preserved 	<ul style="list-style-type: none"> • Q&A • Table
Others	csv	<ul style="list-style-type: none"> • For Q&A documents, each qa pair is treated as one chunk • For table data, each row is considered as a separate chunk 	<ul style="list-style-type: none"> • Context between each pair is preserved 	<ul style="list-style-type: none"> • Q&A • Table

Others	rtf	<ul style="list-style-type: none"> Split the text based on its hierarchical structure (headings and sections), preserving logical divisions and optionally using overlap to maintain context. 	<ul style="list-style-type: none"> Preserves the logical flow of the document. 	Custom
Others	txt	Same as PDF, only extraction will be different		
Others	docx	Same as PDF, only extraction will be different		
Others	json	<ul style="list-style-type: none"> Dividing each text content present in json format Grouping 5-10 messages together to form a chunk 	Easy to process for the LLM and chunk size is consistent.	Custom
Others	xml	<ul style="list-style-type: none"> Split the chat XML data into chunks by grouping a fixed number of<Text> elements (5-10 messages) together. 	Maintains XML structure, easy for LLM to process and chunk size is consistent.	Custom
Others	html	<ul style="list-style-type: none"> Can use HTMLSemanticPreservingSplitter from langchain which splits html content correctly, while keeping semantic meaning intact 	Context is maintained	Custom

7. Knowledge Base

Data Source	Format	Chunking Method	Rationale	Strategy Name
Knowledge Base	.pdf	<ul style="list-style-type: none">Fixed-size chunking with overlap splits where each chunk partially repeats part of the previous one to preserve context	Context is not lost at chunk boundaries, and info flow is maintained.	Custom
Knowledge Base	.txt	<ul style="list-style-type: none">Fixed-size chunking with overlap splits where each chunk partially repeats part of the previous one to preserve context	Context is not lost at chunk boundaries, and info flow is maintained.	Custom

8. Email

Data Source	Format	Chunking Method	Rationale	Strategy Name
Email	.pst	<ul style="list-style-type: none">Extract each email from the .pst file (pypff)Group convo thread or subjectChunk by N consecutive emails in that threadRetain metadata (sender name etc)Overlap of emails	Conversational flow is maintained and overlap ensures no context loss.	Custom

Email	.eml	<ul style="list-style-type: none"> • Extract each email from the .pst file (mailbox) • Group convo thread or subject • Chunk by N consecutive emails in that thread • Retain metadata (sender name etc) • Overlap of emails 	Conversational flow is maintained and overlap ensures no context loss.	Custom
Email	.olm	<ul style="list-style-type: none"> • Extract each email from the .pst file (extract_msg) • Group convo thread or subject • Chunk by N consecutive emails in that thread • Retain metadata (sender name etc) • Overlap of emails 	Conversational flow is maintained and overlap ensures no context loss.	Custom
Email	.mbox	<ul style="list-style-type: none"> • Extract each email from the .pst file (olmreader) • Group convo thread or subject • Chunk by N consecutive emails in that thread • Retain metadata (sender name etc) • Overlap of emails 	Conversational flow is maintained and overlap ensures no context loss.	Custom

9. Packet Capture Format

Data Source	Format	Chunking Method	Rationale	Strategy Name
Packet Capture Format	.pcap	<ul style="list-style-type: none">● Group packets into chunks of 5–10 rows each, of every 5-10 seconds.● Include an overlap of a few packets● Attach the column headers to each for clarity.	Preserves context and is scalable and robust.	Custom