

Practical Lab File for Agentic AI

Sharda School of Engineering and Technology

Name- Kunwar Utkarsh Kant Mishra

Sys ID- 2023438539

Semester/Year- 6th sem, 3rd Year

Faculty-In-Charge/Submitted To

Mr. Ayush Singh



SHARDA
UNIVERSITY
Beyond Boundaries

5 Levels of Text Splitting

for NLP & LLM Pipelines

Abstract

This document explores five hierarchical strategies for text splitting, a critical preprocessing step in Natural Language Processing (NLP) and Large Language Model (LLM) pipelines. Proper text chunking significantly improves embedding quality, retrieval accuracy, and context preservation in downstream tasks. We progress from simple character-based splitting to advanced agentic methods, providing implementation details and best practices for each approach.

Contents

1	Project Overview	3
1.1	Why Text Splitting Matters.....	3
2	Architecture & Processing Flow	3
3	The Five Levels	3
3.1	Overview.....	3
3.2	Installation Requirements.....	4
4	Level 1: Character Splitting	4
4.1	Introduction	4
4.2	Key Concepts	4
4.3	Implementation.....	5
4.3.1	Manual Implementation	5
4.3.2	Using LangChain	5
4.4	Adding Chunk Overlap.....	5
5	Level 2: Recursive Character Splitting	6
5.1	Introduction	6
5.2	Separator Hierarchy.....	6
5.3	Implementation.....	7
5.4	Adding Metadata.....	7
5.5	Optimized Chunking	7
6	Level 3: Document-Specific Splitting	8
6.1	Introduction	8
6.2	Markdown Splitting.....	8
6.2.1	Markdown Separator Hierarchy	8
6.2.2	Implementation.....	9
6.3	Python Code Splitting	9
6.3.1	Python Separator Hierarchy.....	9
6.3.2	Implementation.....	10
6.4	JavaScript Code Splitting	10

6.4.1	JavaScript Separator Hierarchy	10
6.4.2	Implementation.....	11
7	Level 4: Semantic Splitting	11
7.1	Introduction	11
7.2	How It Works	11
7.3	Implementation.....	13
7.4	Understanding Embeddings and Similarity	14
8	Level 5: Agentic Splitting	14
8.1	Introduction	14
8.2	Advantages	14
8.3	Implementation.....	16
8.4	Advanced Agentic Chunking	17
9	Evaluation & Best Practices	17
9.1	Importance of Evaluation	17
9.2	Evaluation Frameworks	18
9.3	Choosing the Right Level.....	18
9.4	Best Practices	18
10	Conclusion	19
10.1	Key Takeaways	19
10.2	Further Resources	19
	Appendix: Quick Reference	19

1 Project Overview

Text splitting, also known as chunking, is a fundamental preprocessing technique that divides large text documents into smaller, manageable pieces. This process is essential for:

Token Limit Management: Working within context window constraints of LLMs

Improved Retrieval: Enhancing semantic search and retrieval accuracy

Better Embeddings: Creating more focused vector representations

Context Preservation: Maintaining semantic coherence within chunks

1.1 Why Text Splitting Matters

Have you ever tried to input a long document into ChatGPT only to receive a "text too long" error? Or struggled to give your application effective long-term memory? Text splitting addresses these challenges by intelligently dividing content while preserving context and meaning.

2 Architecture & Processing Flow

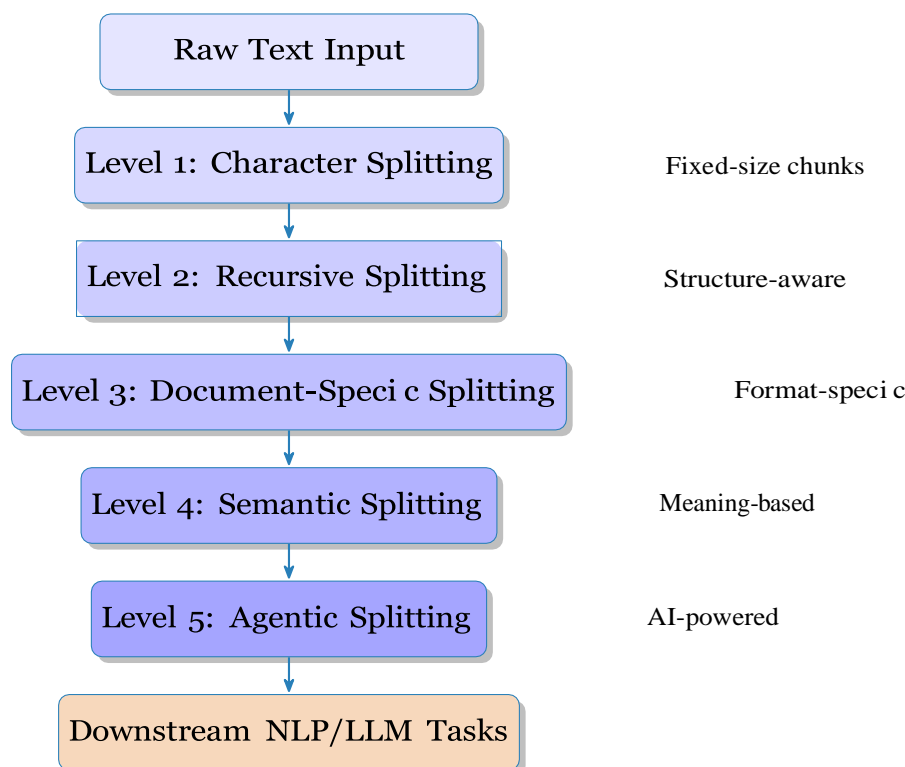


Figure 1: Text Splitting Pipeline Architecture

3 The Five Levels

3.1 Overview

1. **Level 1: Character Splitting** Simple static character chunks

2. **Level 2: Recursive Character Splitting** Recursive chunking based on separators
3. **Level 3: Document-Specific Splitting** Format-aware chunking (PDF, Python, Markdown)
4. **Level 4: Semantic Splitting** Embedding-based chunking
5. **Level 5: Agentic Splitting** LLM-powered intelligent chunking

3.2 Installation Requirements

This tutorial uses code from LangChain. Install the required packages:

Code Cell 0

```
1 pip install langchain langchain-text-splitters
2 pip install langchain_experimental langchain_openai
```

4 Level 1: Character Splitting

4.1 Introduction

Character splitting is the most basic form of text division. It divides text into fixed-size chunks of N characters, regardless of content structure or semantic boundaries.

Advantages & Disadvantages

Pros:

- Simple to implement
- Predictable chunk sizes
- Fast processing

Cons:

- Ignores document structure
- May split words or sentences awkwardly
- Poor context preservation

4.2 Key Concepts

Key Concept

Chunk Size

The number of characters per chunk (e.g., 50, 100, 1000)

Chunk Overlap

The number of characters that overlap between consecutive chunks to preserve context

4.3 Implementation

4.3.1 Manual Implementation

Code Cell 1

```
1 text = "This is the text I would like to chunk up. It is the  
    example text for this exercise"
```

Code Cell 2

```
1 # Create a list to hold chunks  
2 chunks = []  
3 chunk_size = 35 # Characters  
4  
5 # Iterate through text with step size equal to chunk_size  
6 for i in range(0, len(text), chunk_size):  
7     chunk = text[i:i + chunk_size]  
8     chunks.append(chunk)  
9  
10 print(chunks)  
11 # Output: ['This is the text I would like to',  
12 #         'chunk up. It is the example text',  
13 #         ' for this exercise']
```

4.3.2 Using LangChain

In production systems, we work with Document objects that contain both text and metadata for easier manipulation and Itering.

Code Cell 3

```
1 from langchain_text_splitters import CharacterTextSplitter  
2  
3 # Initialize the splitter  
4 text_splitter = CharacterTextSplitter(  
5     chunk_size=35,  
6     chunk_overlap=0,  
7     separator='',  
8     strip_whitespace=False  
9 )  
10  
11 # Create documents  
12 text = "This is the text I would like to chunk up. It is the  
    example text for this exercise"  
13 documents = text_splitter.create_documents([text])  
14  
15 for doc in documents:  
16     print(f"Content: {doc.page_content}")  
17     print(f"Metadata: {doc.metadata}\n")
```

4.4 Adding Chunk Overlap

Chunk overlap helps preserve context by ensuring that adjacent chunks share some content.

Code Cell 4

```

1 # Configure splitter with overlap
2 text_splitter = CharacterTextSplitter (
3     chunk_size=35,
4     chunk_overlap=4,
5     separator=' '
6 )
7
8 documents = text_splitter.create_documents ([text])
9
10 # Notice how chunks now overlap:
11 # Chunk 1: "This is the text I would like to "
12 # Chunk 2: "o chunk up. It is the example text"
13 #      ^^^^ <- overlaps with end of Chunk 1

```

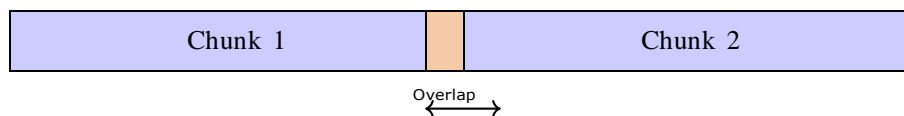


Figure 2: Visualization of chunk overlap

5 Level 2: Recursive Character Splitting

5.1 Introduction

The recursive character text splitter addresses the limitations of Level 1 by respecting document structure. It uses a hierarchy of separators to intelligently split text at natural boundaries.

5.2 Separator Hierarchy

The default separators (in order of priority):

1. "\n\n" Paragraph breaks (double newlines)
2. "\n" Single newlines
3. " " Spaces
4. "" Individual characters

5.3 Implementation

Code Cell 5

```
1 from langchain_text_splitters import
   RecursiveCharacterTextSplitter
2
3 text = """
4 One of the most important things I didn't understand about the
   world
5 when I was a child is the degree to which the returns for
   performance
6 are superlinear.
7
8 Teachers and coaches implicitly told us the returns were linear.
9 "You get out," I heard a thousand times, "what you put in."
10
11 It's obviously true that the returns for performance are
   superlinear
12 in business. Some think this is a flaw of capitalism.
13 """
14
15 # Initialize the splitter
16 text_splitter = RecursiveCharacterTextSplitter (
17     chunk_size=65,
18     chunk_overlap=0
19 )
20
21 documents = text_splitter.create_documents([text])
22
23 for i, doc in enumerate(documents):
24     print(f"Chunk {i+1}: {doc.page_content}\n")
```

5.4 Adding Metadata

Code Cell 6

```
1 # Create documents with metadata
2 docs = text_splitter.create_documents([text])
3
4 for i, doc in enumerate(docs):
5     doc.metadata = {
6         "source_file": "file.txt",
7         "chunk_no": i,
8         "chunk_size": len(doc.page_content)
9     }
10
11 # Access metadata
12 for doc in docs:
13     print(f"Chunk {doc.metadata['chunk_no']}: {doc.metadata}")
```

5.5 Optimized Chunking

The recursive splitter includes intelligent "snapping" to nearest separators:

Code Cell 7

```

1 # Larger chunk size for paragraph-level splits
2 text_splitter = RecursiveCharacterTextSplitter(
3     chunk_size=450,
4     chunk_overlap=0
5 )
6
7 documents = text_splitter.create_documents([text])
8
9 # With chunk_size=450 (or even 469), the splitter creates
10 # perfect paragraph breaks due to separator snapping

```

6 Level 3: Document-Specific Splitting

6.1 Introduction

Different document types require specialized splitting strategies. This level provides format-aware splitters for Markdown, Python, JavaScript, and other formats.

6.2 Markdown Splitting

6.2.1 Markdown Separator Hierarchy

1. `\n#{1,6}` Headers (H1 through H6)
2. ``\n`` Code blocks
3. `\n***+\n` Horizontal lines (asterisks)
4. `\n -+\n` Horizontal lines (dashes)
5. `\n___+\n` Horizontal lines (underscores)
6. `\n\n` Double newlines
7. `\n` Single newlines
8. `" "` Spaces
9. `"""` Characters

6.2.2 Implementation

Code Cell 8

```

1 from langchain_text_splitters import MarkdownTextSplitter
2
3 markdown_text = """
4 # Fun in California
5
6 ## Driving
7 Try driving on the 1 down to San Diego
8
9 ### Food
10 Make sure to eat a burrito while you're there
11
12 ## Hiking
13 Go to Yosemite
14 """
15
16 splitter = MarkdownTextSplitter(
17     chunk_size=40,
18     chunk_overlap=0
19 )
20
21 documents = splitter.create_documents([markdown_text])
22
23 # Chunks will align with markdown sections
24 for doc in documents:
25     print(f"Chunk: {doc.page_content}\n")

```

6.3 Python Code Splitting

6.3.1 Python Separator Hierarchy

1. \n\nclass Class definitions
2. \n\ndef Function definitions
3. \n\tdef Indented functions (methods)
4. \n\n Double newlines
5. \n Single newlines
6. " " Spaces
7. """ Characters

6.3.2 Implementation

Code Cell 9

```

1 from langchain_text_splitters import PythonCodeTextSplitter
2
3 python_text = """
4 class Person:
5     def __init__(self, name, age):
6         self.name = name
7         self.age = age
8
9 p1 = Person("John", 36)
10
11 for i in range(10):
12     print(i)
13 """
14
15 python_splitter = PythonCodeTextSplitter(
16     chunk_size=100,
17     chunk_overlap=0
18 )
19
20 documents = python_splitter.create_documents([python_text])
21
22 # Chunks preserve class and function boundaries
23 for i, doc in enumerate(documents):
24     print(f"Chunk {i+1}:\n{doc.page_content}\n")

```

6.4 JavaScript Code Splitting

6.4.1 JavaScript Separator Hierarchy

Separator	Purpose
\nfunction	Function declarations
\nconst	Constant variable declarations
\nlet	Block-scoped variable declarations
\nvar	Variable declarations
\nclass	Class definitions
\nif	Conditional statements
\nfor	For loops
\nwhile	While loops
\nswitch	Switch statements
\ncase	Case clauses in switch
\ndefault	Default clause in switch
\n\n	Double newlines
\n	Single newlines
" "	Spaces
""	Characters

Table 1: JavaScript separator hierarchy

6.4.2 Implementation

Code Cell 10

```
1 from langchain_text_splitters import
   RecursiveCharacterTextSplitter, Language
2
3 javascript_text = """
4 // Function is called, the return value will end up in x
5 let x = myFunction(4, 3);
6
7 function myFunction(a, b) {
8     // Function returns the product of a and b
9     return a * b;
10 }
11 """
12
13 js_splitter = RecursiveCharacterTextSplitter.from_language(
14     language=Language.JS,
15     chunk_size=65,
16     chunk_overlap=0
17 )
18
19 documents = js_splitter.create_documents([javascript_text])
20
21 for doc in documents:
22     print(f"Chunk: {doc.page_content}\n")
```

7 Level 4: Semantic Splitting

7.1 Introduction

Semantic chunking uses embeddings to group text based on meaning rather than structure. This approach identifies natural topic boundaries by analyzing semantic similarity between sentences.

7.2 How It Works

1. Split text into sentences
2. Generate embeddings for each sentence
3. Calculate similarity between consecutive sentences
4. Create chunk boundaries where similarity drops below a threshold

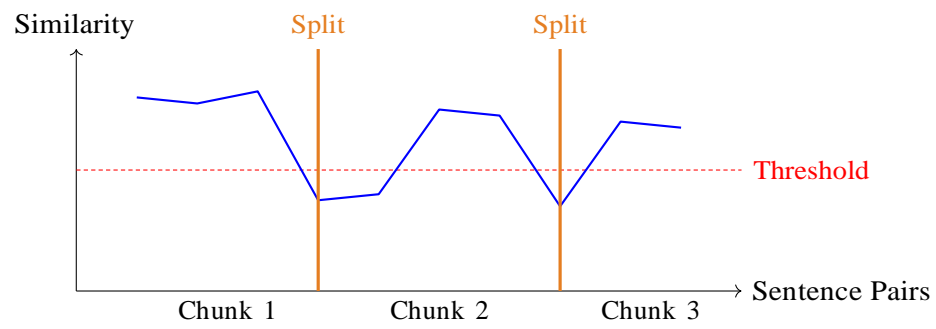


Figure 3: Semantic chunking based on similarity threshold

7.3 Implementation

Code Cell 11

```

1 from langchain_experimental.text_splitter import SemanticChunker
2 from langchain_openai import OpenAIEmbeddings
3 import os
4
5 # Set up API key
6 os.environ['OPENAI_API_KEY'] = 'your-api-key-here'
7
8 # Sample text with distinct topics
9 tesla_text = """Tesla's Q3 Results
10 Tesla reported record revenue of $25.2B in Q3 2024.
11 The company exceeded analyst expectations by 15%.
12 Revenue growth was driven by strong vehicle deliveries.
13
14 Model Y Performance
15 The Model Y became the best-selling vehicle globally,
16 with 350,000 units sold.
17 Customer satisfaction ratings reached an all-time high of 96%.
18 Model Y now represents 60% of Tesla's total vehicle sales.
19
20 Production Challenges
21 Supply chain issues caused a 12% increase in production costs.
22 Tesla is working to diversify its supplier base.
23 New manufacturing techniques are being implemented to reduce
    costs."""
24
25 # Initialize semantic chunker
26 semantic_splitter = SemanticChunker(
27     embeddings=OpenAIEmbeddings(),
28     breakpoint_threshold_type="percentile",
29     breakpoint_threshold_amount=70
30 )
31
32 # Split text semantically
33 chunks = semantic_splitter.split_text(tesla_text)
34
35 print("SEMANTIC CHUNKING RESULTS:")
36 print("=" * 50)
37 for i, chunk in enumerate(chunks, 1):
38     print(f"Chunk {i}: ({len(chunk)} chars)")
39     print(f'"{chunk}"\n')

```

7.4 Understanding Embeddings and Similarity

Code Cell 12

```

1 from openai import OpenAI
2 import numpy as np
3 from sklearn.metrics.pairwise import cosine_similarity
4
5 client = OpenAI()
6
7 sentences = [
8     "Tesla reported record revenue of $25.2B in Q3 2024.",
9     "The company exceeded analyst expectations by 15%",
10    "Revenue growth was driven by strong vehicle deliveries.",
11    "The Model Y became the best-selling vehicle globally.",
12    "Customer satisfaction ratings reached an all-time high of 96%."
13 ]
14
15 # Generate embeddings
16 embeddings = []
17 for s in sentences:
18     response = client.embeddings.create(
19         model="text-embedding-3-small",
20         input=s
21     )
22     embeddings.append(response.data[0].embedding)
23
24 embeddings = np.array(embeddings)
25
26 # Calculate similarity between consecutive sentences
27 similarity_matrix = cosine_similarity(embeddings)
28
29 print(f"Similarity S1 vs S2: {similarity_matrix[0][1]:.3f}")
30 print(f"Similarity S2 vs S3: {similarity_matrix[1][2]:.3f}")
31 print(f"Similarity S3 vs S4: {similarity_matrix[2][3]:.3f}")
32 print(f"Similarity S4 vs S5: {similarity_matrix[3][4]:.3f}")
33
34 # High similarity (>0.8) = keep together
35 # Low similarity (<0.6) = create boundary

```

8 Level 5: Agentic Splitting

8.1 Introduction

Agentic splitting leverages LLMs to intelligently determine chunk boundaries. The AI analyzes content, identifies logical topics, and creates chunks that preserve semantic coherence.

8.2 Advantages

Context-aware: Understands content meaning and structure

Flexible: Adapts to different document types automatically

Intent-driven: Can optimize for specific use cases

Natural boundaries: Creates semantically coherent chunks

8.3 Implementation

Code Cell 13

```

1 from langchain_openai import ChatOpenAI
2
3 # Initialize the LLM
4 llm = ChatOpenAI(model="gpt-4", temperature=0)
5
6 # Text to chunk
7 tesla_text = """Tesla's Q3 Results
8 Tesla reported record revenue of $25.2B in Q3 2024.
9 The company exceeded analyst expectations by 15%.
10 Revenue growth was driven by strong vehicle deliveries.
11
12 Model Y Performance
13 The Model Y became the best-selling vehicle globally,
14 with 350,000 units sold.
15 Customer satisfaction ratings reached an all-time high of 96%.
16 Model Y now represents 60% of Tesla's total vehicle sales.
17
18 Production Challenges
19 Supply chain issues caused a 12% increase in production costs.
20 Tesla is working to diversify its supplier base.
21 New manufacturing techniques are being implemented to reduce
    costs."""
22
23 # Create the prompt
24 prompt = f"""
25 You are a text chunking expert. Split this text into logical
    chunks.
26
27 Rules:
28 - Each chunk should be around 200 characters or less
29 - Split at natural topic boundaries
30 - Keep related information together
31 - Put "<<<SPLIT>>>" between chunks
32
33 Text:
34 {tesla_text}
35
36 Return the text with <<<SPLIT>>> markers where you want to split:
37 """
38
39 # Get AI response
40 print("Asking AI to chunk the text...")
41 response = llm.invoke(prompt)
42 marked_text = response.content
43 print(marked_text)
44
45 # Split the text at the markers
46 chunks = marked_text.split("<<<SPLIT>>>")
47
48 # Clean up the chunks
49 clean_chunks = [chunk.strip() for chunk in chunks if chunk.strip()]
50
51 # Show results
52 print("\nAGENTIC CHUNKING RESULT_16:")
53 print("=" * 50)
54 for i, chunk in enumerate(clean_chunks, 1):
55     print(f"Chunk {i}: ({len(chunk)} chars)")
56     print(f'"{chunk}"\n')

```

8.4 Advanced Agentic Chunking

Code Cell 14

```
1 # More sophisticated prompt for domain-specific chunking
2 advanced_prompt = f"""
3 You are an expert at chunking technical documents for RAG systems
4 .
5 Context: This text will be used in a question-answering system
6 about Tesla's financial performance.
7
8 Requirements:
9 1. Each chunk must be self-contained and understandable
10 2. Target chunk size: 150-250 characters
11 3. Preserve numerical data with context
12 4. Keep section headers with their content
13 5. Mark splits with <<<SPLIT>>>
14
15 Optimization goals:
16 - Maximize retrieval accuracy
17 - Preserve causal relationships
18 - Maintain temporal context
19
20 Text:
21 {tesla_text}
22
23 Provide the chunked text:
24 """
25
26 # Process with advanced prompt
27 response = llm.invoke(advanced_prompt)
28 # ... (same processing as before)
```

9 Evaluation & Best Practices

9.1 Importance of Evaluation

Chunking strategy significantly impacts downstream performance. Always evaluate your chunks using retrieval metrics:

Retrieval Accuracy: Are the right chunks retrieved?

Context Completeness: Do chunks contain sufficient context?

Answer Quality: Does chunking improve final outputs?

9.2 Evaluation Frameworks

Framework	Description
LangChain Evals	Comprehensive evaluation tools for RAG systems https://python.langchain.com/docs/guides/evaluation/
Llama Index Evals	Specialized retrieval evaluation metrics https://docs.llamaindex.ai/en/stable/module_guides/evaluating/
RAGAS	Advanced RAG assessment framework https://github.com/explodinggradients/ragas

Table 2: Evaluation frameworks for text chunking

9.3 Choosing the Right Level

Level	Best For	Avoid For
Character	Simple prototypes, uniform text	Production systems, varied content
Recursive	General-purpose text, mixed content	Code, structured documents
Document-Specific	Code files, Markdown, specific formats	General prose
Semantic	Topic-based retrieval, coherent sections	Performance-critical systems
Agentic	Complex documents, custom requirements	High-volume processing

Table 3: Chunking strategy selection guide

9.4 Best Practices

1. **Start Simple:** Begin with recursive splitting before advanced methods
2. **Measure Impact:** Always benchmark against retrieval performance
3. **Consider Cost:** Semantic and agentic methods require API calls
4. **Use Overlap:** 10-20% overlap helps preserve context
5. **Test Multiple Sizes:** Optimal chunk size varies by use case (typically 200-1000 chars)
6. **Add Metadata:** Include source, chunk number, and context in metadata
7. **Preserve Structure:** Maintain headings, sections, and formatting
8. **Iterate:** Chunking is not one-size-fits-all; experiment and refine

10 Conclusion

Text splitting is a critical preprocessing step that directly impacts the performance of NLP and LLM applications. By understanding and applying these five levels from simple character splitting to sophisticated agentic methods you can significantly improve embedding quality, retrieval accuracy, and overall system performance.

10.1 Key Takeaways

- Choose the splitting strategy that matches your document type and use case
- Always evaluate chunking decisions against downstream task performance
- Consider the trade-offs between simplicity, accuracy, and computational cost
- Use overlap to preserve context across chunk boundaries
- Iterate and refine based on real-world performance metrics

10.2 Further Resources

LangChain Text Splitters Documentation: https://python.langchain.com/docs/modules/data_connection/document_transformers/

ChunkViz Visualization Tool: <https://www.chunkviz.com>

OpenAI Embeddings Guide: <https://platform.openai.com/docs/guides/embeddings>

Appendix: Quick Reference

Chunk Size Guidelines

Use Case	Recommended Size	Overlap
Question Answering	200-500 chars	50-100 chars
Summarization	500-1000 chars	100-200 chars
Code Analysis	100-300 chars	20-50 chars
Long-form Retrieval	1000-2000 chars	200-400 chars

Table 4: Chunk size recommendations by use case

Common Pitfalls

1. **Too Small:** Chunks lack sufficient context
2. **Too Large:** Exceeds model context windows, poor retrieval precision
3. **No Overlap:** Important context lost at boundaries
4. **Excessive Overlap:** Increased storage and processing costs
5. **Ignoring Structure:** Character splitting on formatted documents
6. **No Evaluation:** Optimizing without measuring impact