# Loading and Processing Data

## Instructor

Dipanjan Sarkar

Head of Community & Principal AI Scientist at Analytics Vidhya

Google Developer Expert - ML & Cloud Champion Innovator

Published Author

# Key Objectives of the Project

Build a Search Engine on

Wikipedia Articles & Research Papers

Process text and PDF documents

Create document and contextual chunks

Index chunks and embeddings in a vector database

Experiment with different retrievers

Loading and Processing Data

Analytics Vidhya

# Data Sources



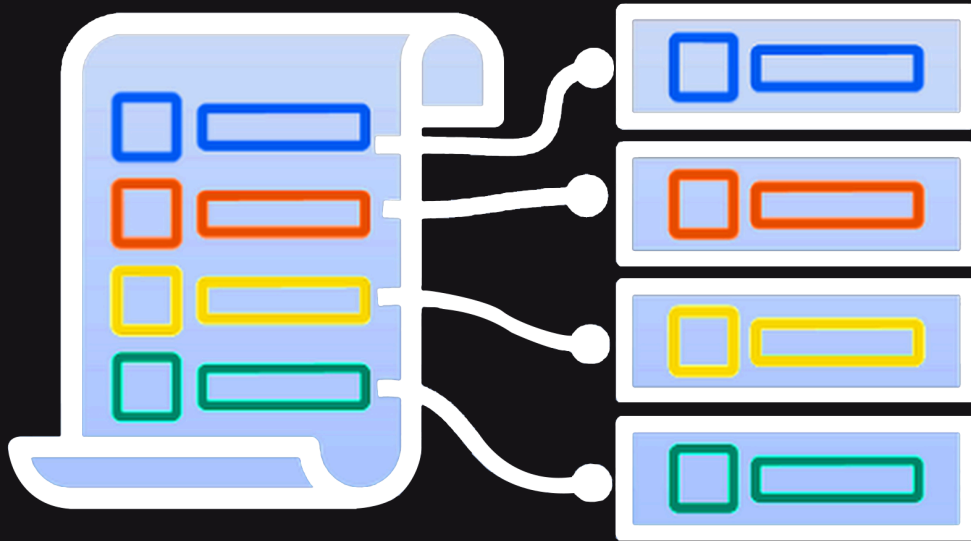Text Article JSON
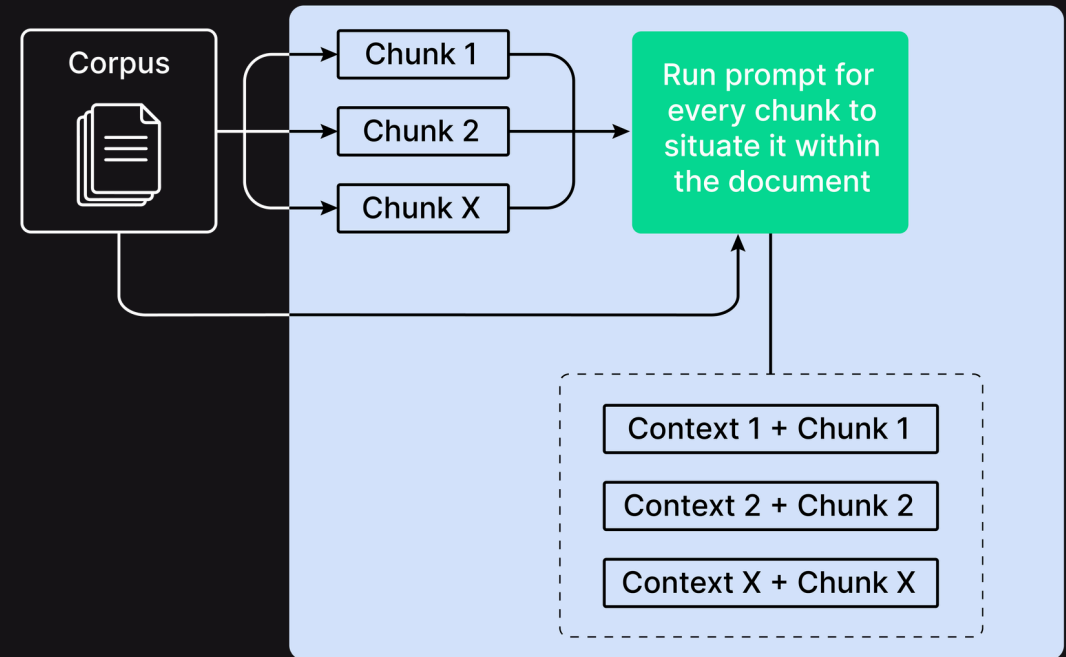


Research Paper PDFs

# Data Loader



JSON Loader

PDF Loader

# Chunking Strategies



Recursive Character
Text Splitting & Chunking

Contextual Chunking

# Thank You

Analytics
Vidhya