

# Node Parsers

## Instructors

Prashant Sahu

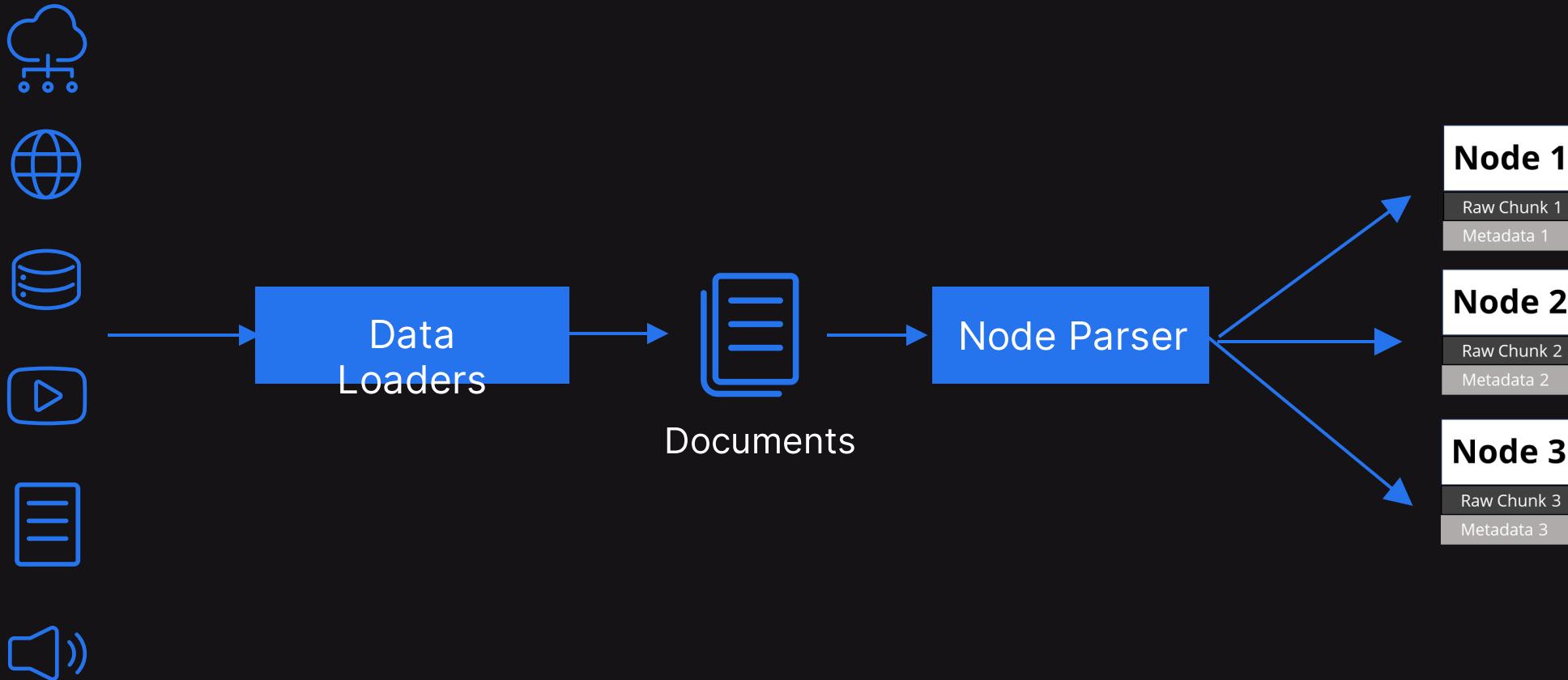
Manager - Data Science, Analytics Vidhya

Ravi Theja

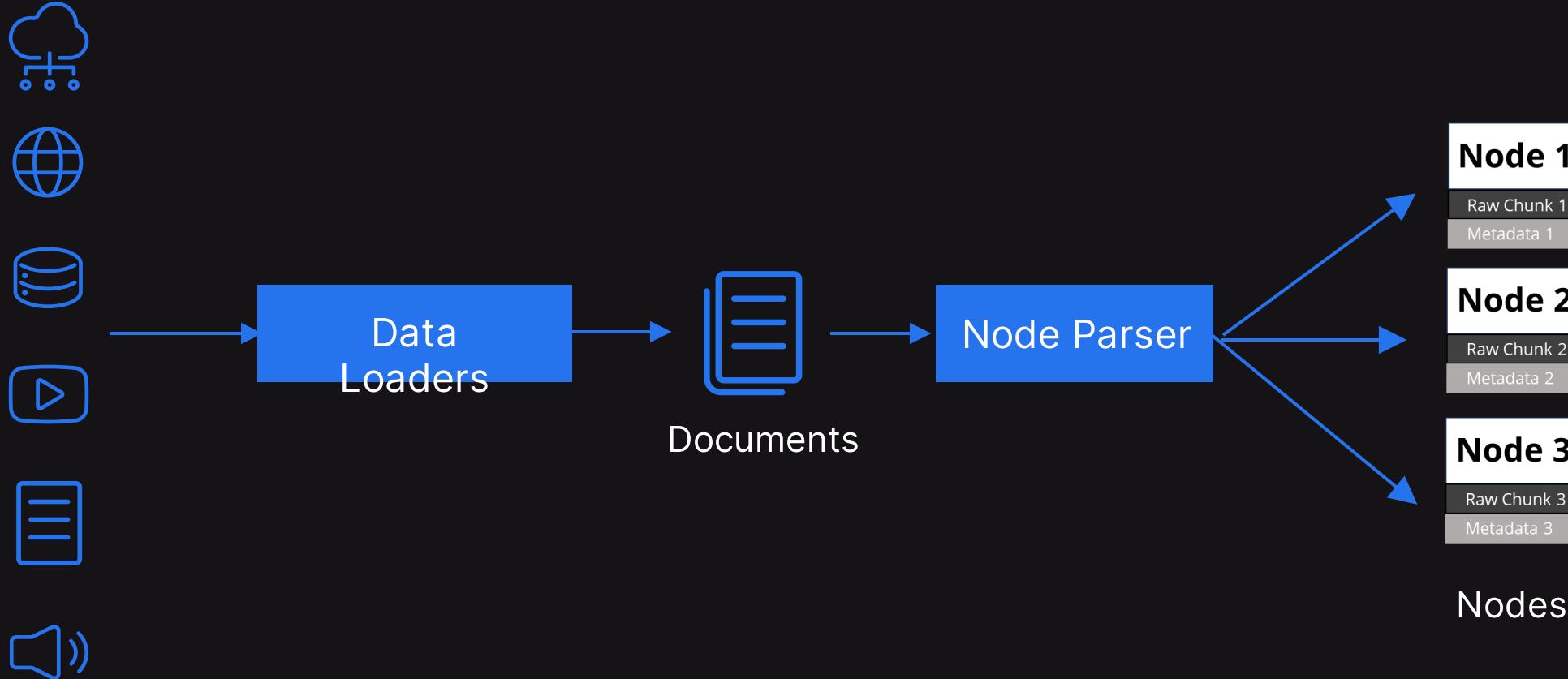
Developer Advocate Engineer, LlamalIndex



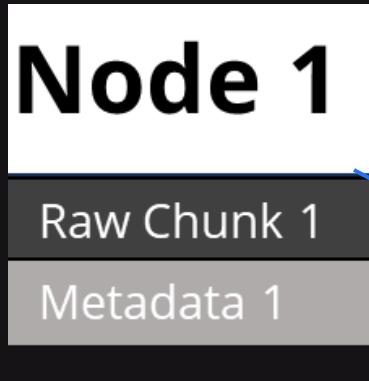
# Documents in Llamalndex



# Nodes in Llamalndex



# Node



```
nodes[0].text
```

'In this part, we\nwill discuss the corresponding configurations for four major\nparts of the Transformer, includin\ng normalization, position\nembeddings, activation functions, and attention and bias.\nTo make this survey more self\n-contained, we present the\ndetailed formulations for these configurations in Table 6.\nNormalization Methods. Trai\nning instability is a challeng-\ning issue for pre-training LLMs. To alleviate this issue,\nnormalization is a wide\nly adopted strategy to stabilize the\\ntraining of neural networks. In the vanilla Transformer [22],\nLayerNorm [25\n6] is employed. Recently, several advanced\\nnormalization techniques have been proposed as alterna-\ntives to Layer\nNorm, e.g., RMSNorm, and DeepNorm.\n•LayerNorm. In the early research, BatchNorm [265] is\\na commonly used normaliz\nation method. However, it is\\ndifficult to deal with sequence data of variable lengths and\\nsmall-batch data. Thus,\nLayerNorm [256] is introduced to\\nconduct layerwise normalization. Specifically, the mean and\\nvariance over all ac

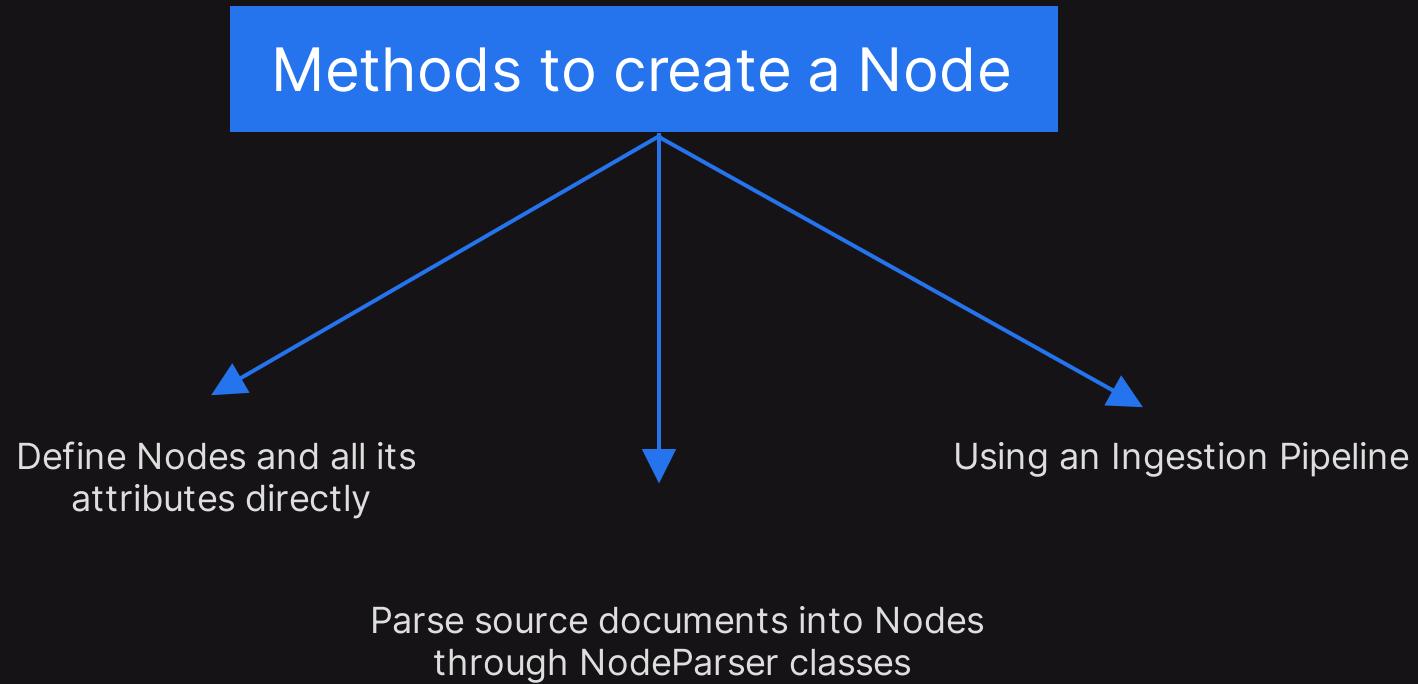
```
nodes[0].metadata
```

{'page\_label': '23',  
 'file\_name': 'A Survey of Large Language Models.pdf',  
 'file\_path': 'A Survey of Large Language Models.pdf',  
 'file\_type': 'application/pdf',  
 'file\_size': 5664979,  
 'creation\_date': '2024-05-14',  
 'last\_modified\_date': '2024-05-14'}

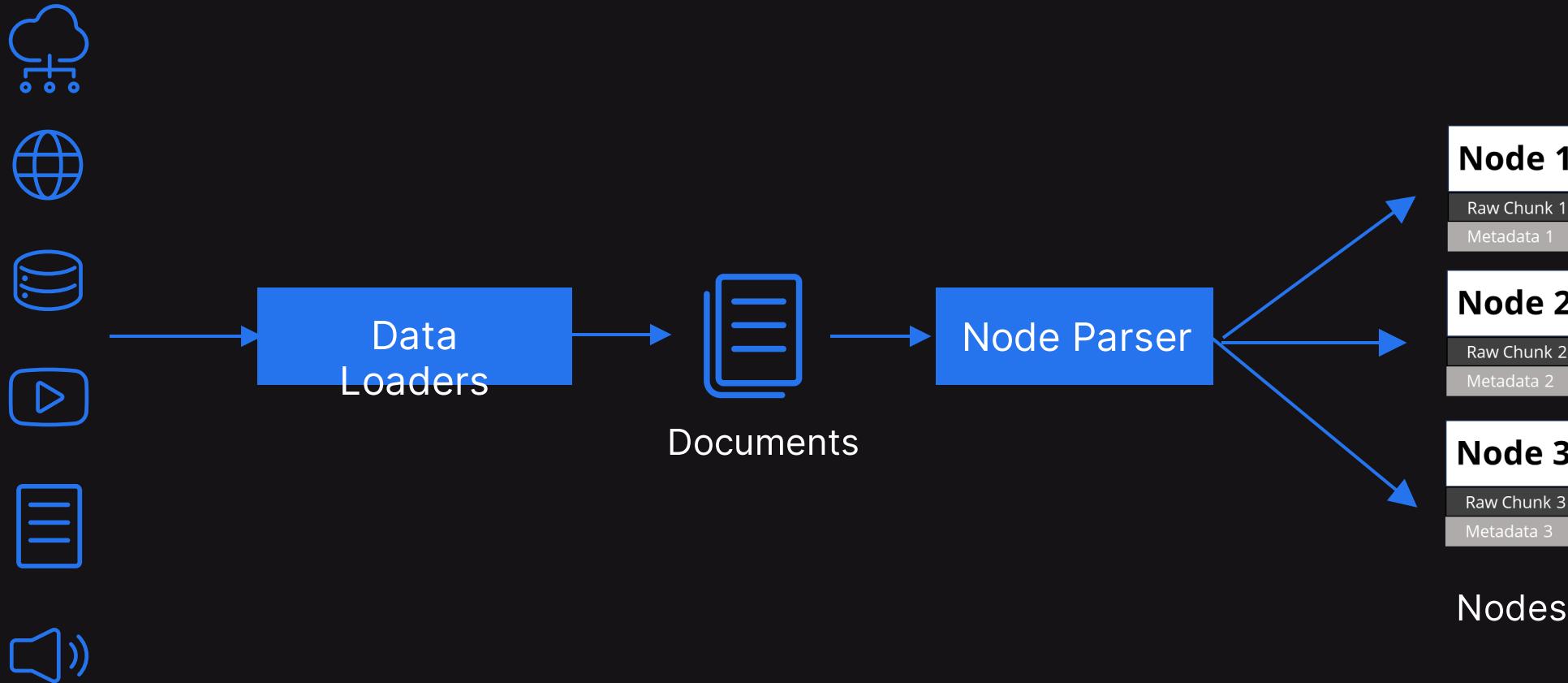
Raw Chunk for Node

Metadata for Node

# Methods to create Nodes



# Introduction to Node Parsers



# Node Parsers

**Node Parsers** break down data into manageable pieces (nodes) based on specific rules or formats, such as text, HTML, JSON, or CSV.

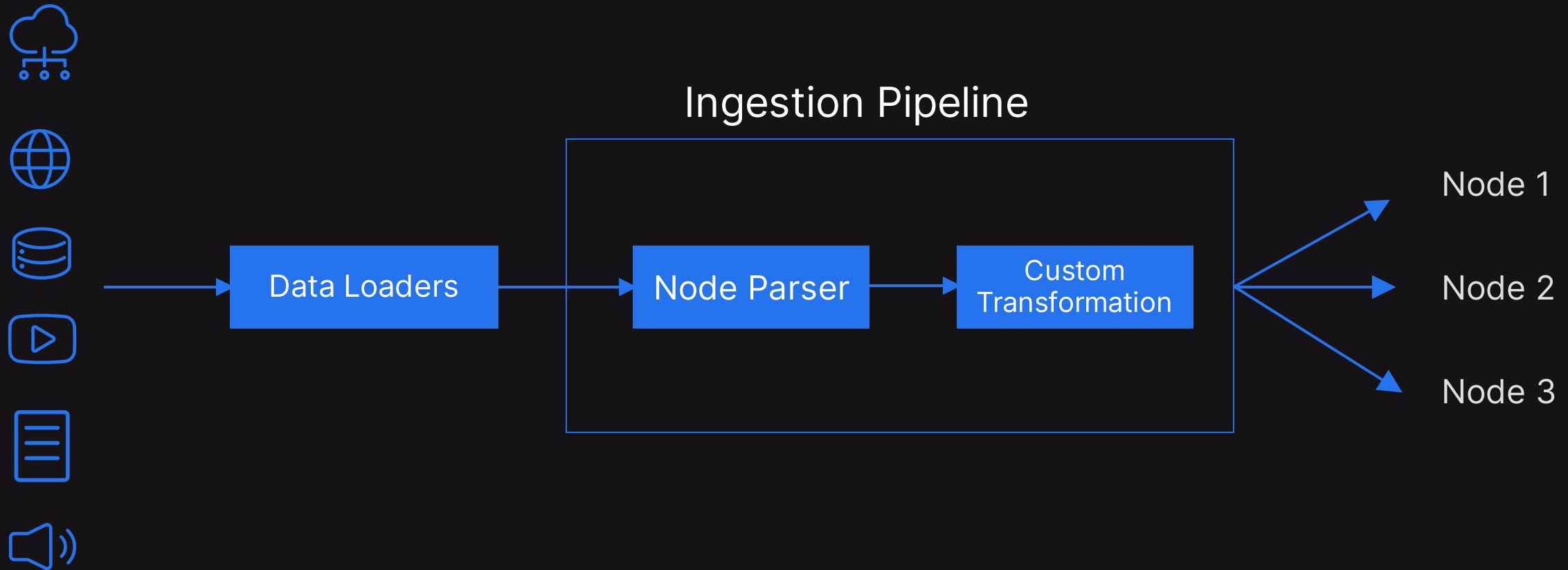
# Node Parsers in Llama Index

- **TextNodeParser**  
For parsing plain text from documents and articles
- **HTMLNodeParser**  
For parsing HTML documents
- **JSONNodeParser**  
Processes JSON data structures
- **CSVNodeParser**  
Handles CSV files by extracting data from rows

# Node Parsers in Llama Index

- **MarkdownNodeParser**  
Specifically designed for Markdown files
- **HierarchicalNodeParser**  
Creates a recursive hierarchy of nodes by splitting documents into multiple levels of nodes.
- **SimpleFileNodeParser**  
Handling of different file formats within a single application.
- **SimpleNodeParser**  
Splits documents into text nodes using a sentence splitter by default.

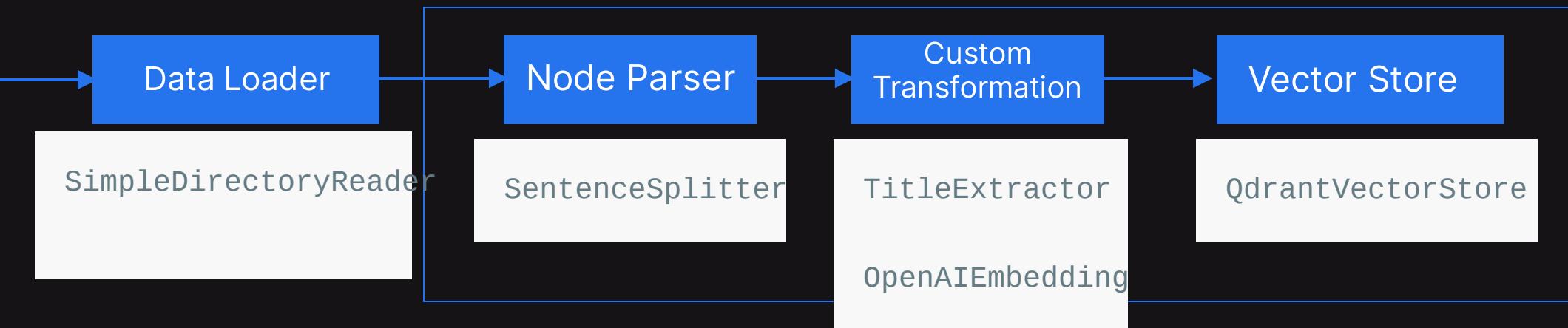
# Ingestion Pipeline



# Ingestion Pipeline: Example



```
pipeline = IngestionPipeline(  
    transformations=[  
        SentenceSplitter(chunk_size=500, chunk_overlap=25),  
        TitleExtractor(),  
        OpenAIEmbedding(),  
    ],  
    vector_store=QdrantVectorStore(),  
)
```



# Thank You