# Semantic Plankton: Multimodal LLMs and RAG for Automated Ocean Microorganism Analysis

## AUTOMATIC IMAGE ANALYSIS WITH PLANKTOSCOPE IN THE PLANDYO PROJECT

Jaronchai Dilokkalayakul [1], Akane Kitamura [2], Takeshi Obayashi [1,2] ( [1]Graduate School of Information Sciences, Tohoku University | [2]WPI-AIMEC )

## Summary

Traditional plankton identification methods are labor-intensive and hard to scale, posing challenges for monitoring marine ecosystems and detecting environmental changes. While CNNs automate classification, they depend on large labeled datasets and lack contextual reasoning. To address this, we propose a framework that combines large language models (LLMs) with retrieval-augmented generation (RAG) to classify plankton using both image and contextual metadata. By retrieving semantically similar examples from a curated vector database and integrating them into the Large Language Model's input, our system enables context-aware, scalable classification with reduced reliance on labeled data.
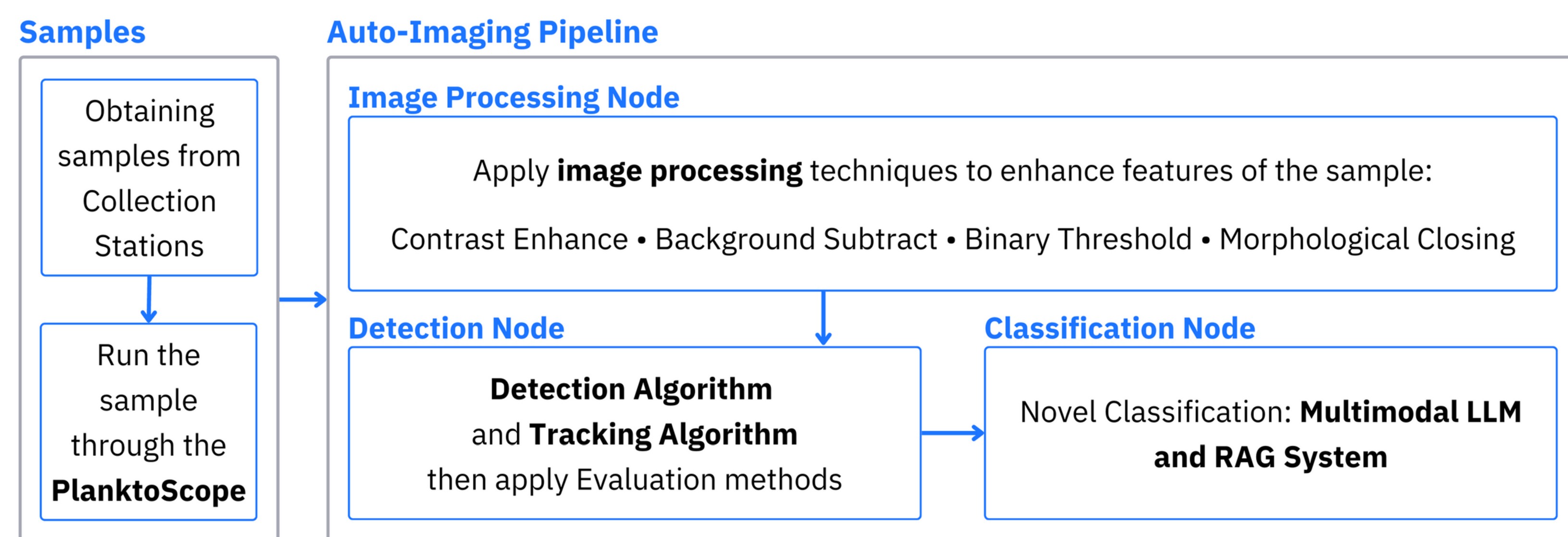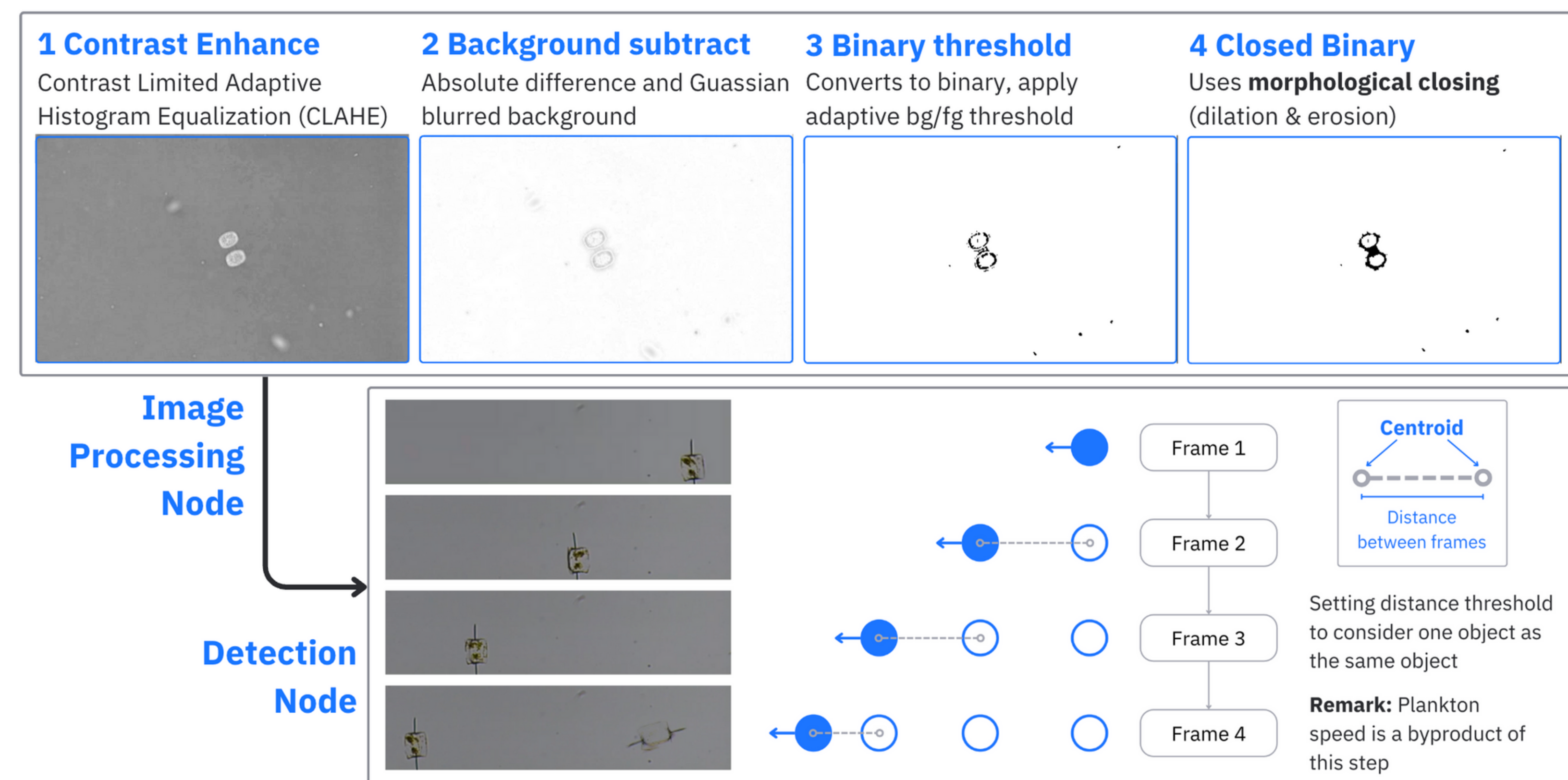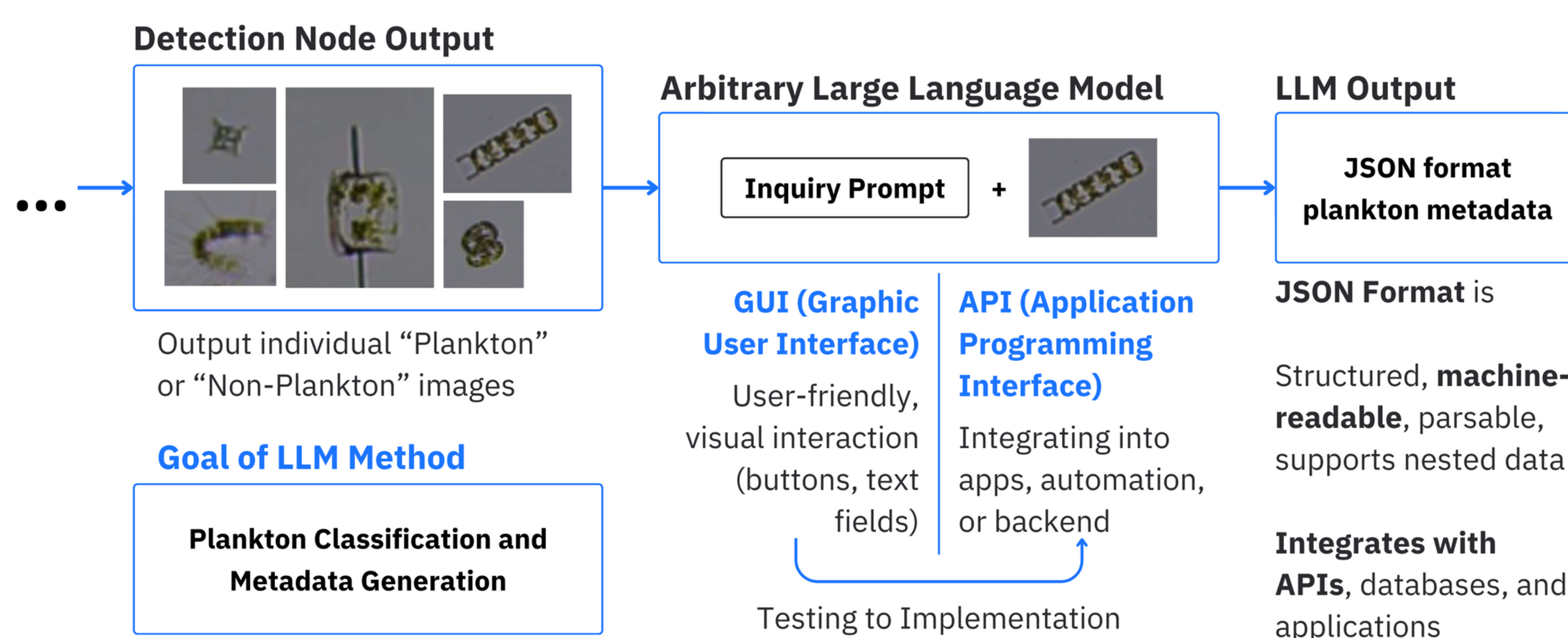
## Process Pipeline

**Samples**

Obtaining samples from Collection Stations

Run the sample through the **PlanktoScope**

**Auto-Imaging Pipeline**

**Image Processing Node**

Apply **image processing** techniques to enhance features of the sample:

Contrast Enhance • Background Subtract • Binary Threshold • Morphological Closing

**Detection Node**

**Detection Algorithm** and **Tracking Algorithm** then apply Evaluation methods

**Classification Node**

Novel Classification: **Multimodal LLM and RAG System**

## Image Processing

**1 Contrast Enhance**
Contrast Limited Adaptive Histogram Equalization (CLAHE)

**2 Background subtract**
Absolute difference and Guassian blurred background

**3 Binary threshold**
Converts to binary, apply adaptive bg/fg threshold

**4 Closed Binary**
Uses **morphological closing** (dilation & erosion)

**Image Processing Node**

**Detection Node**

Frame 1
Frame 2
Frame 3
Frame 4

**Centroid**

Distance between frames

Setting distance threshold to consider one object as the same object

**Remark:** Plankton speed is a byproduct of this step

## Key Concept of LLM Inquiry

**Detection Node Output**

Output individual "Plankton" or "Non-Plankton" images

**Goal of LLM Method**

Plankton Classification and Metadata Generation

**Arbitrary Large Language Model**

Inquiry Prompt +

**GUI (Graphic User Interface)**
User-friendly, visual interaction (buttons, text fields)

**API (Application Programming Interface)**
Integrating into apps, automation, or backend

Testing to Implementation

**LLM Output**

JSON format plankton metadata

JSON Format is

Structured, **machine-readable**, parsable, supports nested data

**Integrates with APIs**, databases, and applications

## Advantage of LLM Implementation

| Feature | Traditional Image Classification | Large Language Model + RAG System |
|---|---|---|
| **Input Data** | Static images | Continuous video frames |
| **Output Format** | Static segmentation | Context-aware, richer taxonomic outputs |
| **Motion Analysis** | ✕ Not possible | Tracks movement trajectories |
| **Environmental Factors** | ✕ Not considered | Includes temperature, pH, angle |
| **Behavioral Insights** | ✕ Limited to morphology | Behavior changes based on factors |

## Example Output Metadata Results



```
{
  "family": "Chaetocerotaceae",
  "genus": "Chaetoceros"
}
```

```
{
  "family": "Ditylaceae",
  "genus": "Ditylum"
}
```

```
{
  "family": "Ceratiaceae",
  "genus": "Ceratium"
}
```

## Large Language Model and Retrieval-Augmented Generation Framework

### RAG implemented LLM with Feedback Loop

**Existing RAG implemented LLM**

**Classification Inquiry Node**

**Our Input**

Prompt + Picture

all-mpnet-base-v2

Embed Model

Embedded vectors

$$\begin{bmatrix} 1 & 2 & n \\ a_{i1} & a_{i2} & \dots & a_{in} \end{bmatrix}$$

Context Injection

**Arbitrary LLM Model**

$$cos\theta = \frac{\sum_0^{n-1}(a_i \cdot b_i)}{\sqrt{\sum_0^{n-1} a_i^2} \cdot \sqrt{\sum_0^{n-1} b_i^2}}$$

Vector Distance

**JSON & Metadata**

```
{
  "family": "Ceratiaceae",
  "genus": "Ceratium",
  "species": "Ceratium furca"
}
```

LLM Output

**Vector Embedding Node**

Vector Database

Append new knowledge into the knowledge-base

Embedded vectors

$$\begin{matrix} & 1 & 2 & n \\ 1 & a_{11} & a_{12} & a_{1n} \\ 2 & a_{21} & a_{22} & a_{2n} \\ 3 & a_{31} & a_{32} & a_{3n} \\ m & a_{m1} & a_{m2} & a_{mn} \end{matrix}$$

Embedded vectors

Picture/PDF/Text of Plankton data

Embed Model

**Embeddings**

$$\begin{matrix} & 1 & 2 & n \\ 1 & a_{11} & a_{12} & a_{1n} \\ 2 & a_{21} & a_{22} & a_{2n} \\ 3 & a_{31} & a_{32} & a_{3n} \\ m & a_{m1} & a_{m2} & a_{mn} \end{matrix}$$

Embed Model

New data from LLM result

**Result Verification with Expert in the field**

New fact-checked and verified data

A prompt and an image of a plankton specimen are provided as the input query. Then the image and textual context are encoded into vector embeddings. The embedding is used to search a vector database containing embedded plankton documents, this retrieves semantically similar examples. The retrieved context is then injected into an LLM, which uses both the input and retrieved information to generate a structured JSON taxonomic metadata output. New outputs are verified by experts to confirm the result. Validated examples are embedded and added back to the vector database, as a continuous improvement to the system's contextual grounding.

## Contact Information

Presenter: **Jaronchai Dilokkalayakul**
Email: dilokkalayakul.jaronchai.p8@dc.tohoku.ac.jp
Personal Website: https://jaronchai.com