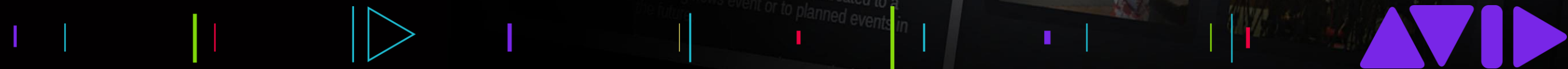


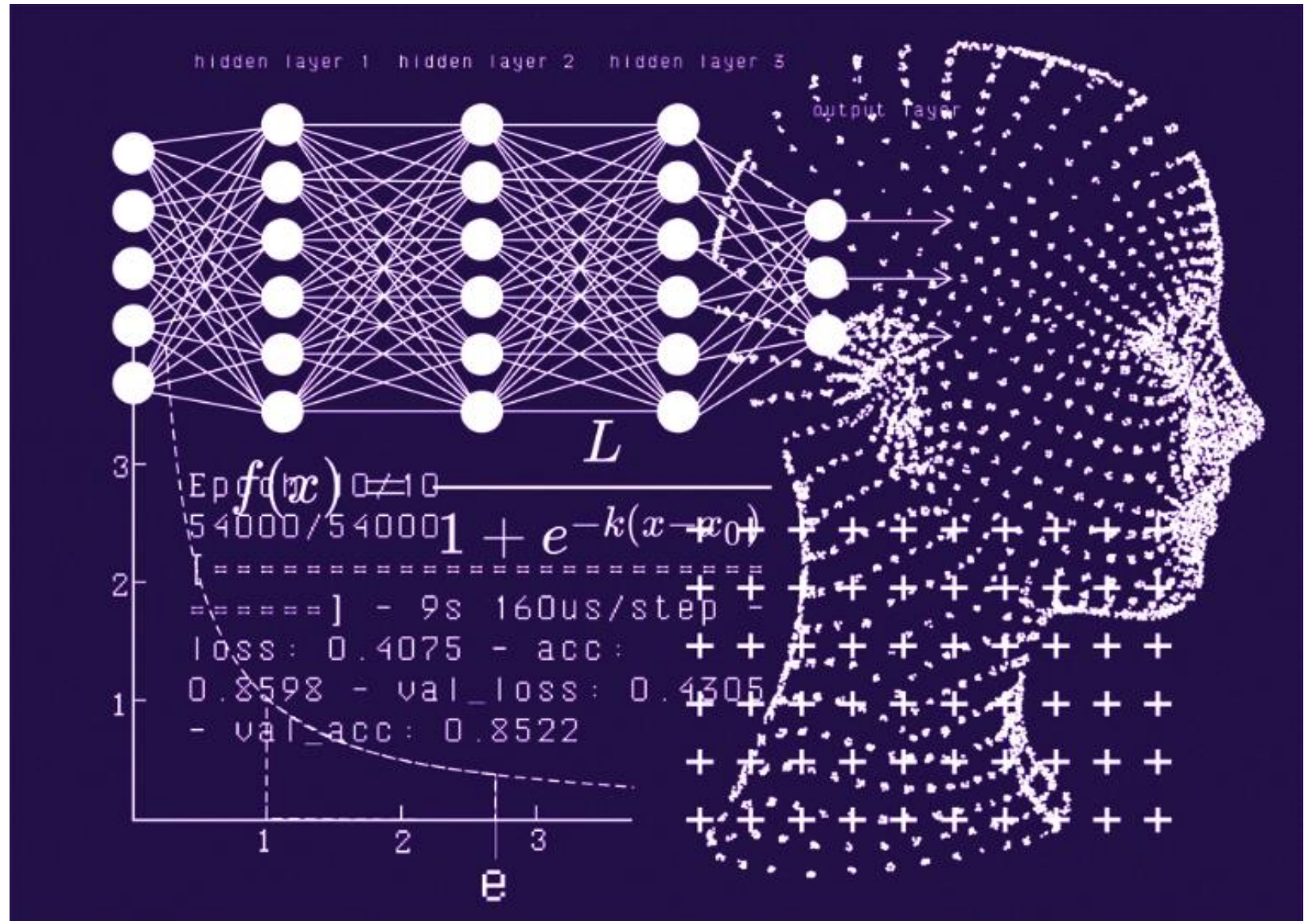
# USE CASE: TRANSFORMING MEDIA ASSET MANAGEMENT WITH AI-DRIVEN SEMANTIC CONTENT DISCOVERY

Rob Gonsalves, Avid  
May 29, 2024



# TRANSFORMING MEDIA ASSET MANAGEMENT WITH AI-DRIVEN SEMANTIC CONTENT DISCOVERY

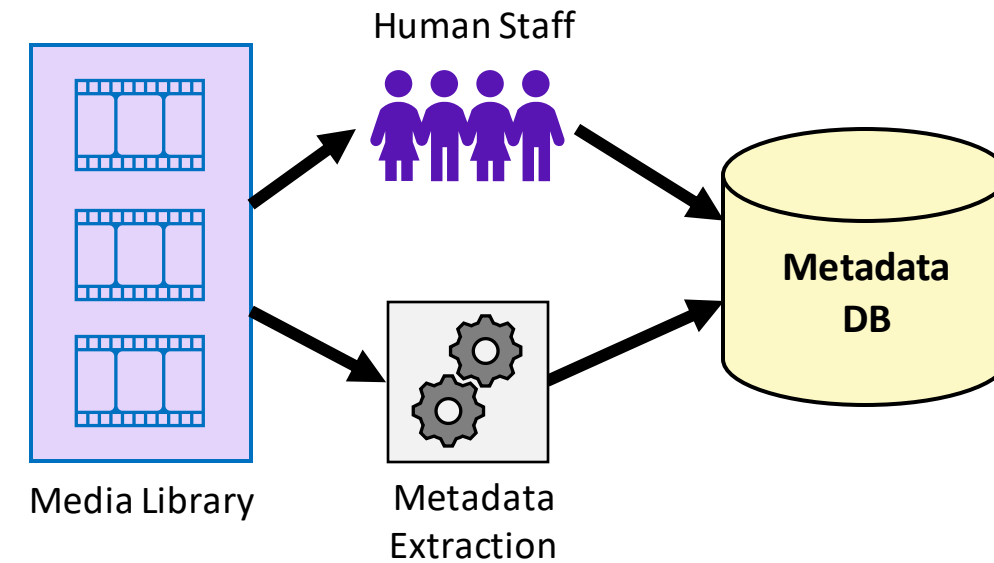
- Introduction to Semantic Content Discovery
- Planning and Implementation
- Outcomes and Evaluations
- Challenges, Successes, and Next Steps
- Q&A



# TRADITIONAL SEARCH METHODS VS. AI/ML APPROACHES

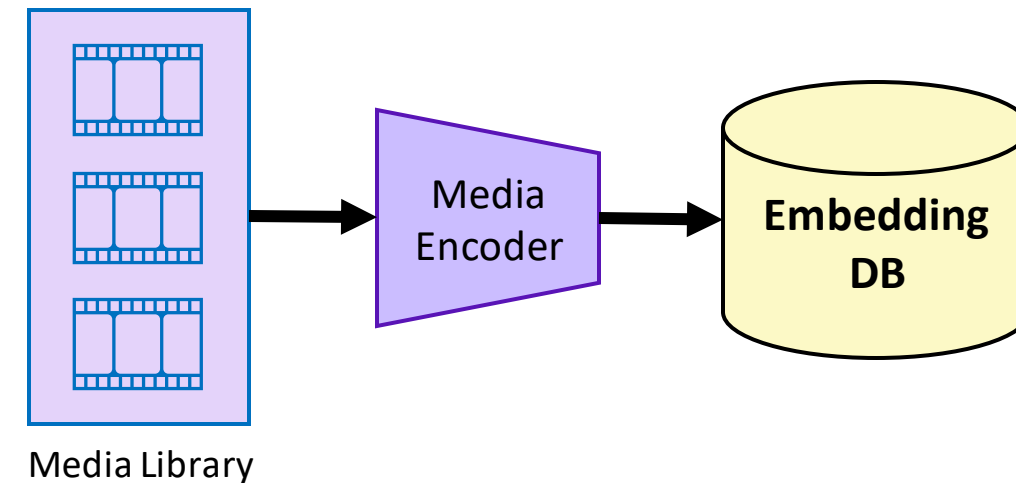
## Metadata-based content discovery:

1. Extract metadata from media
  - Done by humans and/or automated processes
2. Organize into taxonomy (optional)
3. Search using keywords or known terms
4. The system returns media that “hits”



## AI-based semantic content discovery:

1. Create embeddings from media
  - Automated process
  - Contains semantic understanding
2. Search with keywords or phrases
3. The system returns closest matching media



# PLANNING AND IMPLEMENTATION

## Objectives:

- Implement Semantic Content Discovery of Video Media
- Multilingual Capabilities
- Free and Open-source models

## Methods:

- Compare and Evaluate Open-source Semantic Search Models
- Use Metrics like Cosine Similarity and Recall@1 for Analysis

## Implementation:

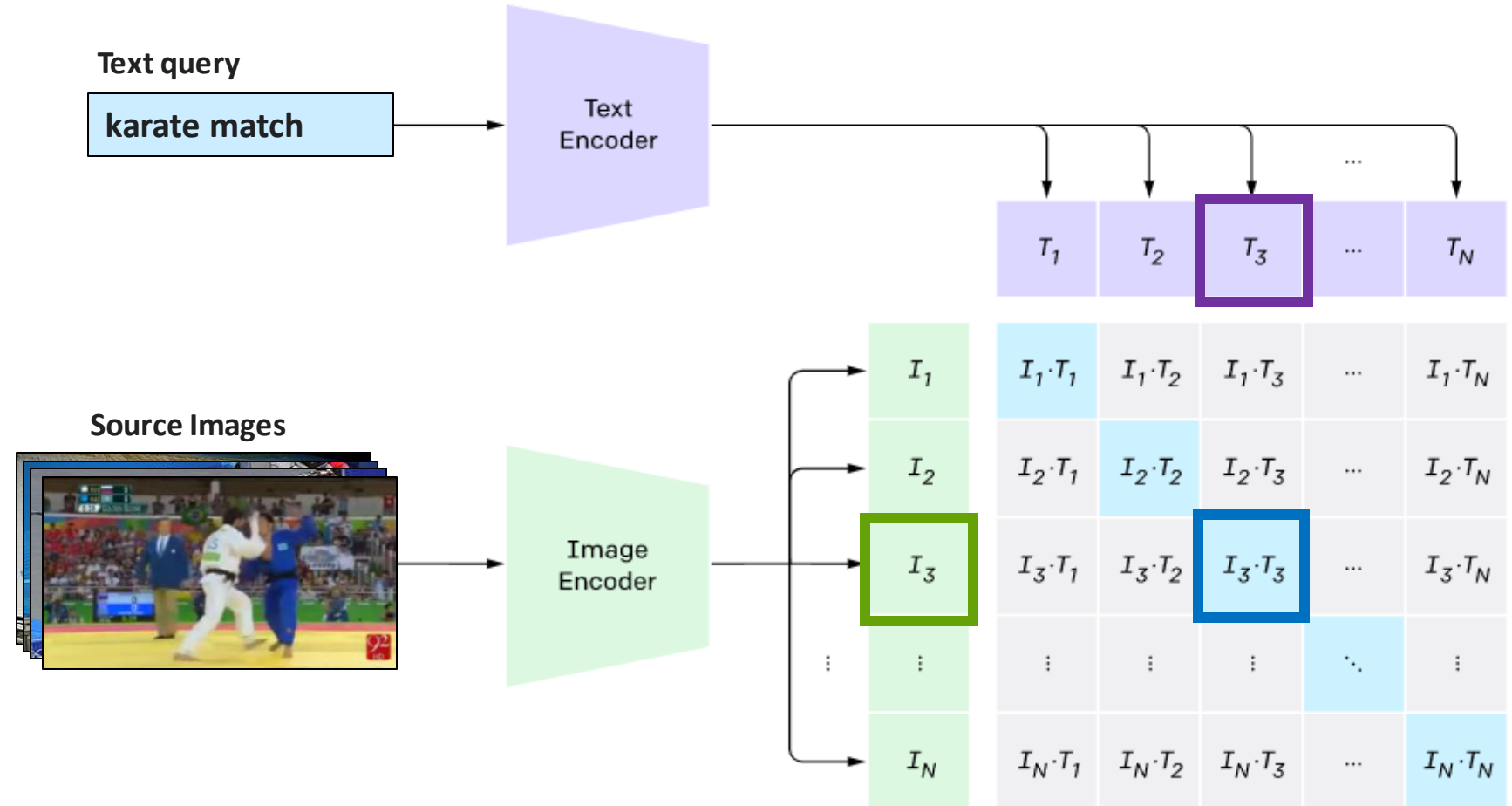
- Prototype in Python
- Productize in C++



# CLIP - IMAGE AND TEXT ENCODERS

## CONTRASTIVE LANGUAGE-IMAGE PRE-TRAINING

- OpenAI created the CLIP models for multi-modal semantic search
- The Text Encoder converts a phrase into an embedding, a list of 512 floating-point numbers
- The Image Encoder converts an image into a similar embedding
- We can use these models to find matching images using an unstructured text query
- English Only

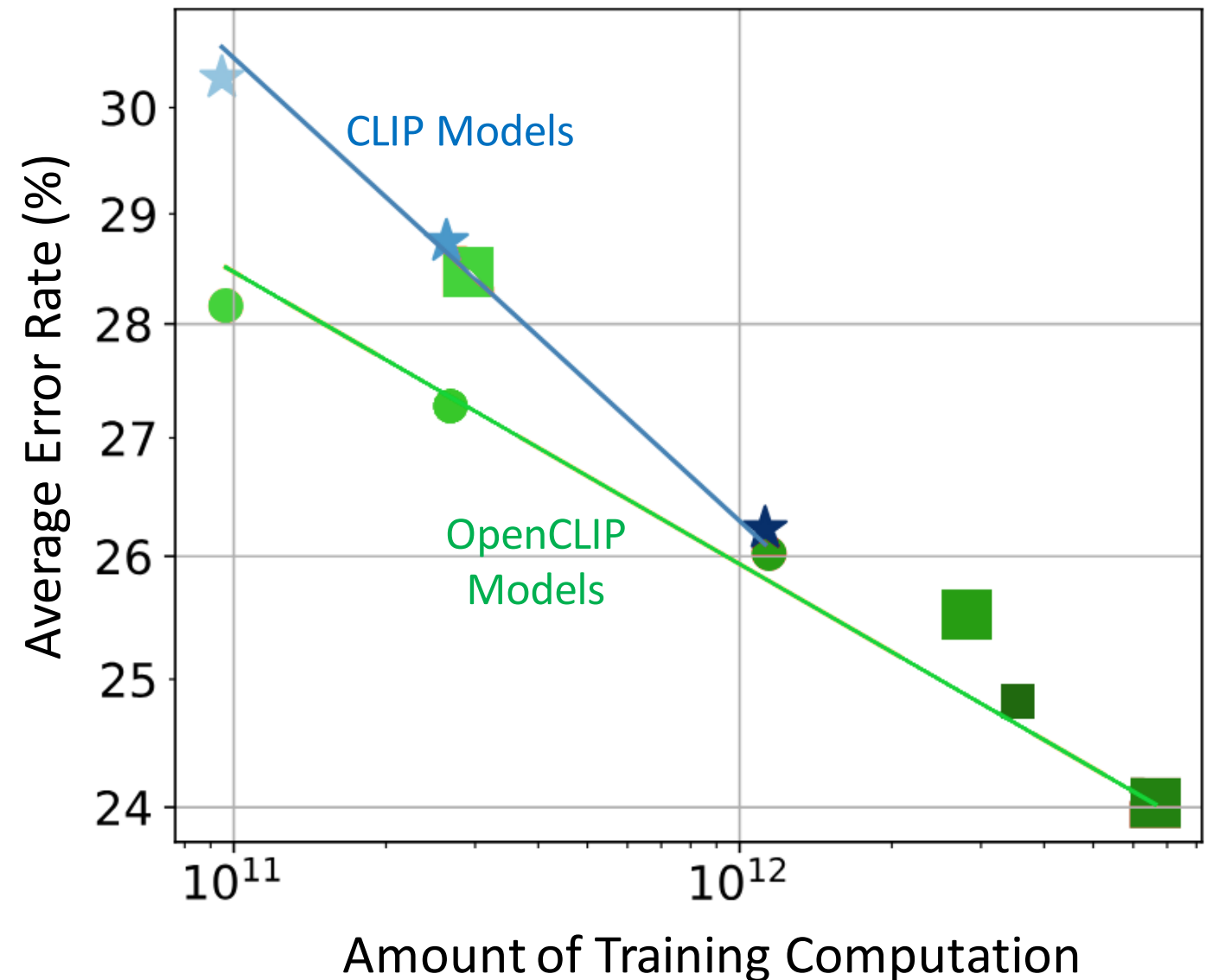


# OpenCLIP from LAION

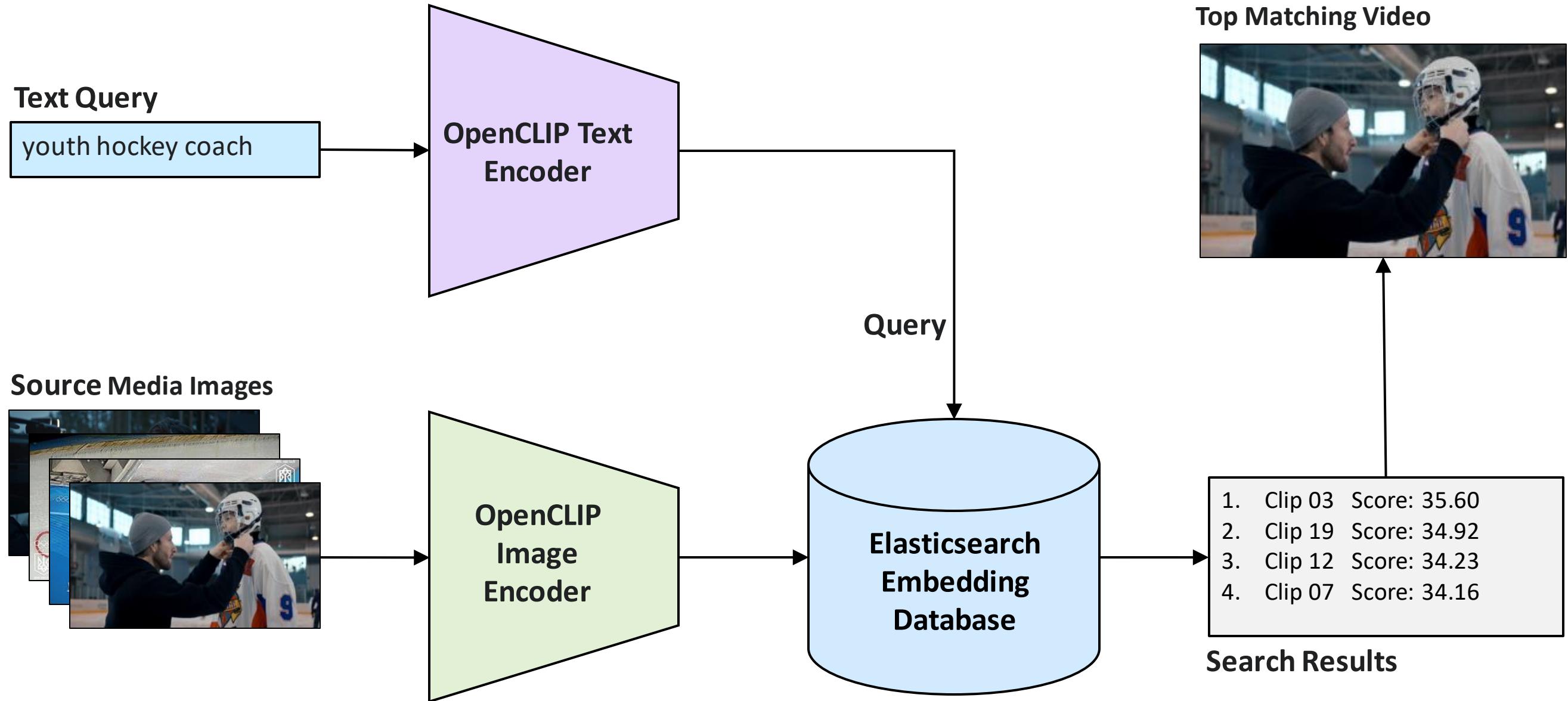
- OpenCLIP is open-source implementation of OpenAI's CLIP.
- The OpenCLIP models were trained on up to two billion image-text pairs.
- Different scaling behaviors were observed between OpenAI and OpenCLIP models despite similar architectures and training recipes.
- Evaluated xlm-roberta-base-ViT-B-32 model using checkpoint laion5b\_s13b\_b90k
- The evaluation workflows, datasets, and models were open-sourced.

## Comparison of Error Rates Between CLIP and OpenCLIP Models

(smaller is better)



# SEMANTIC TEXT AND IMAGE ENCODERS



AI

B I U Aa ab A [List Icons] [Clipboard] [Undo] [Redo] [Print] [Share] [AI]

Claudia Roth is the presenter of MediaCentral Moments

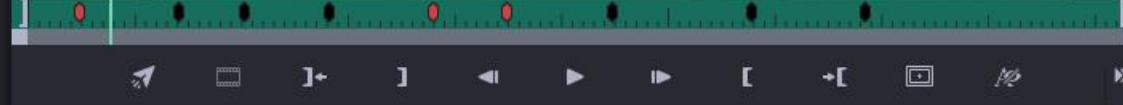
### Avid Ada Recommendations

results: 6

**A-Roll** B-Roll

Score: 81.90%	Score: 80.70%
Score: 78.30%	Score: 77.40%
Score: 76.80%	Score: 76.20%

01:00:00:00 [s] 01:00:00:00 [01:01:42:11]



Au... H... R, A... Meta... File... Storyb... Grap... Traco... Trac...

L R L R

S M S M

Vol 0 A1 Vol 0 A2

Vol 0 Master

Stereo

+0 -4 -8 -14 -20 -26 -30 -35 -40 -46 -52 -58 -60



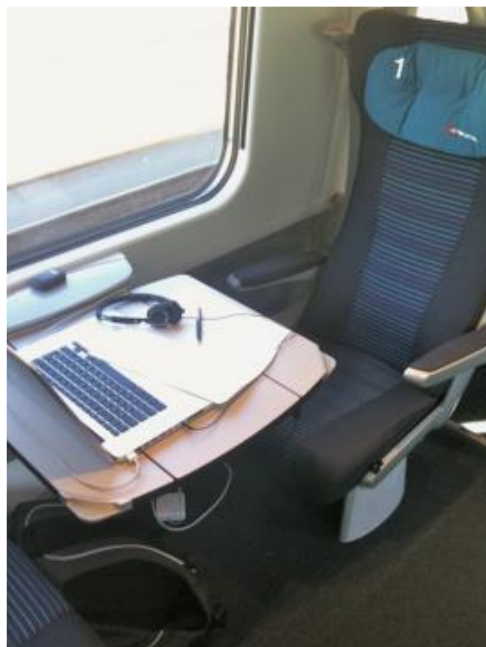
# OUTCOMES AND EVALUATIONS

- Search results from CLIP and OpenCLIP
- Used Google's Crossmodal-3600 dataset for testing
- Compare results for English and multilingual searches for both models with Cosine Similarity and Recall@1 metrics
- Compare image encoding times for both models



## Results from CLIP, search term: "A laptop on a train."

Score: 0.3582



Score: 0.2998



Score: 0.2908



Score: 0.2690

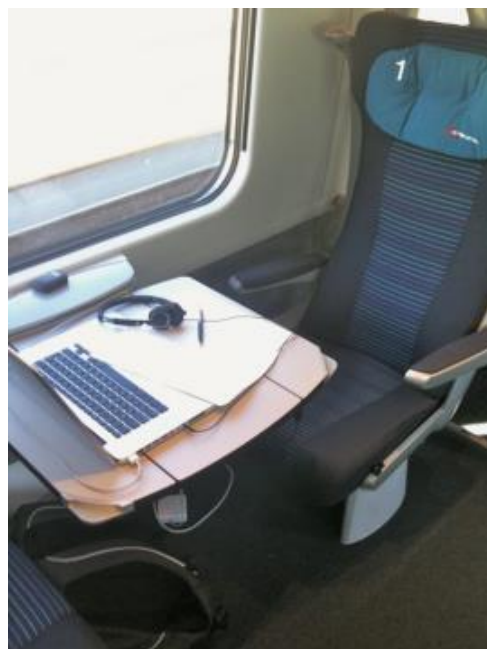


Score: 0.2673



## Results from OpenCLIP, search term: "A laptop on a train."

Score: 0.3669



Score: 0.3069



Score: 0.2964



Score: 0.2778



Score: 0.2659



# Results from CLIP, search term: "Two people talking in a living room."

Score: 0.2927



Score: 0.2920



Score: 0.2808



Score: 0.2681



Score: 0.2644



# Results from OpenCLIP, search term: "Two people talking in a living room."

Score: 0.2656



Score: 0.1886



Score: 0.1874



Score: 0.1860



Score: 0.1743



# CLIP - English

Median of matches (diagonal): 0.309

Median of misses (off-diagonal): 0.149

Median Match-to-Miss Delta: 0.160



A young man smiling and posing.

A group of people having fun in rafting on the water.

A macro shot of a candy floss cover.

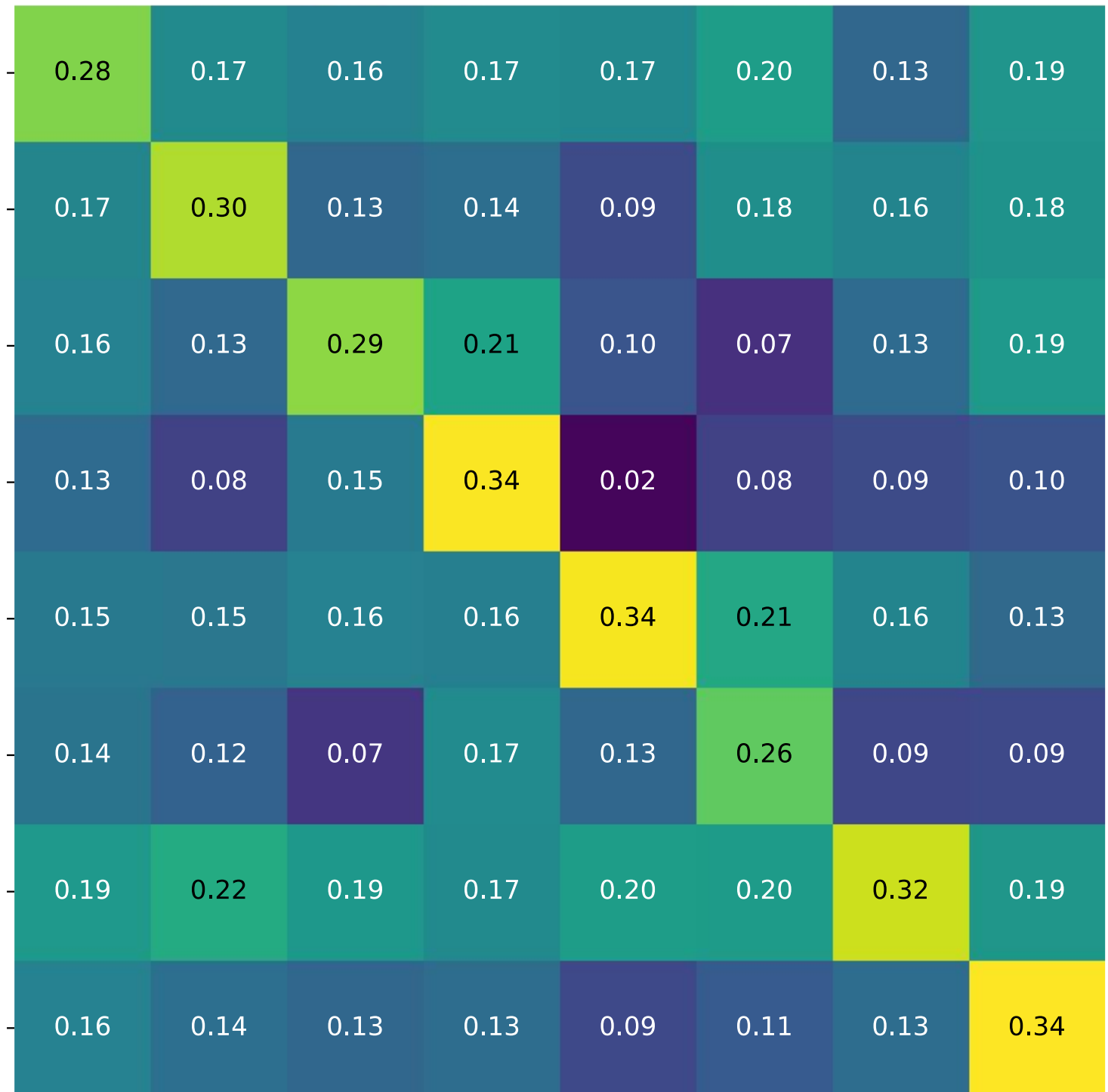
A macro shot of white and pink flowers with a background of white clouds in a blue sky.

An interior view of an old abandoned factory.

An old house in the village through the trees.

The log sheet.

Chinese steamed dim sum in various colors placed in a basket.



# OpenCLIP - English

Median of matches (diagonal): 0.249

Median of misses (off-diagonal): 0.013

Median Match-to-Miss Delta: 0.236



A young man smiling and posing.

A group of people having fun in rafting on the water.

A macro shot of a candy floss cover.

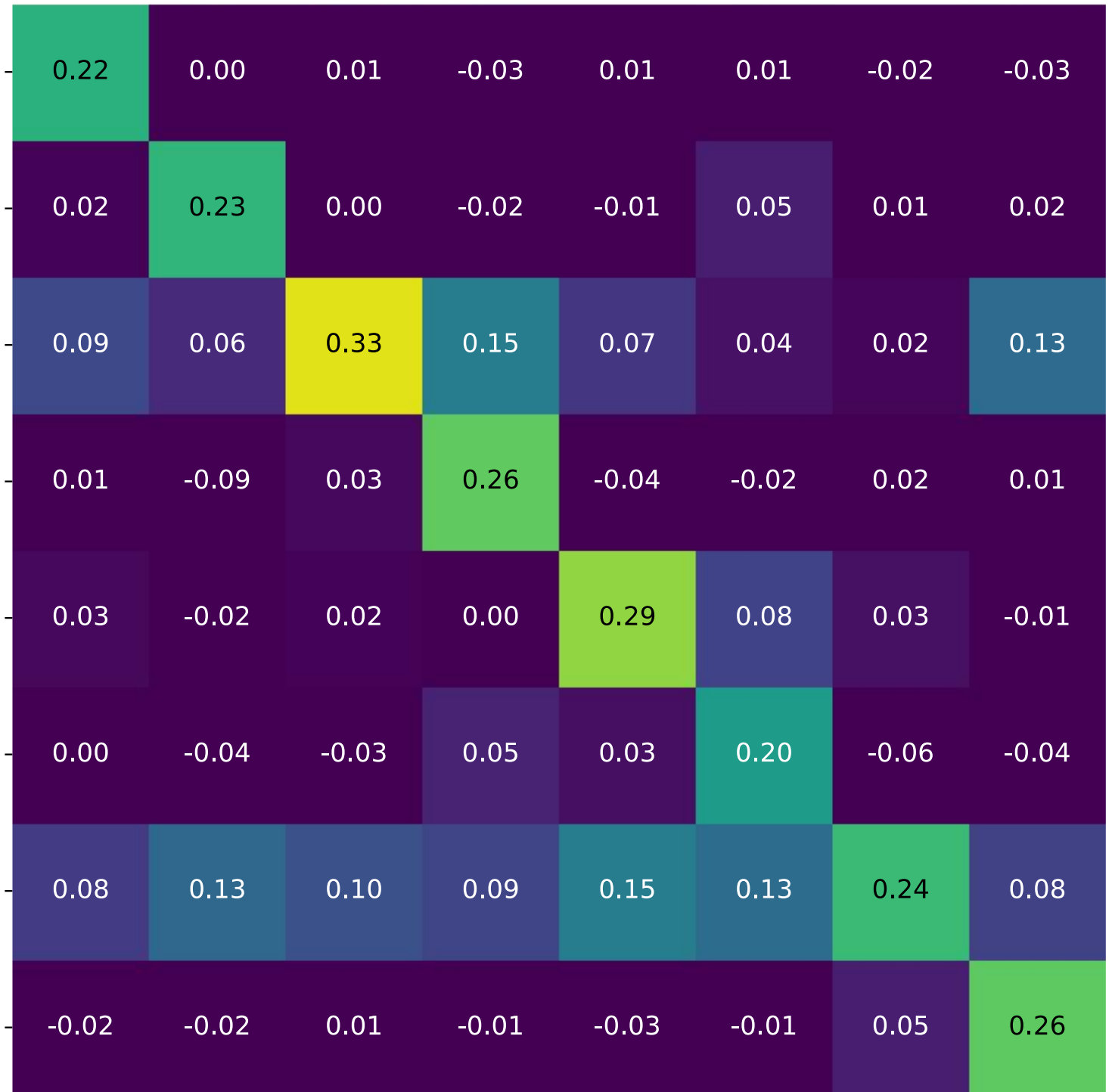
A macro shot of white and pink flowers with a background of white clouds in a blue sky.

An interior view of an old abandoned factory.

An old house in the village through the trees.

The log sheet.

Chinese steamed dim sum in various colors placed in a basket.



# CLIP - Multilingual

Median of matches (diagonal): 0.227

Median of misses (off-diagonal): 0.175

Median Match-to-Miss Delta: 0.052



(English) A young man smiling and posing.

(French) 7 personnes sur une boué jaune et orange faisant du rafting en vertical dans l'eau

(Korean) 타투가 들어있는 사탕 케이스 껍질

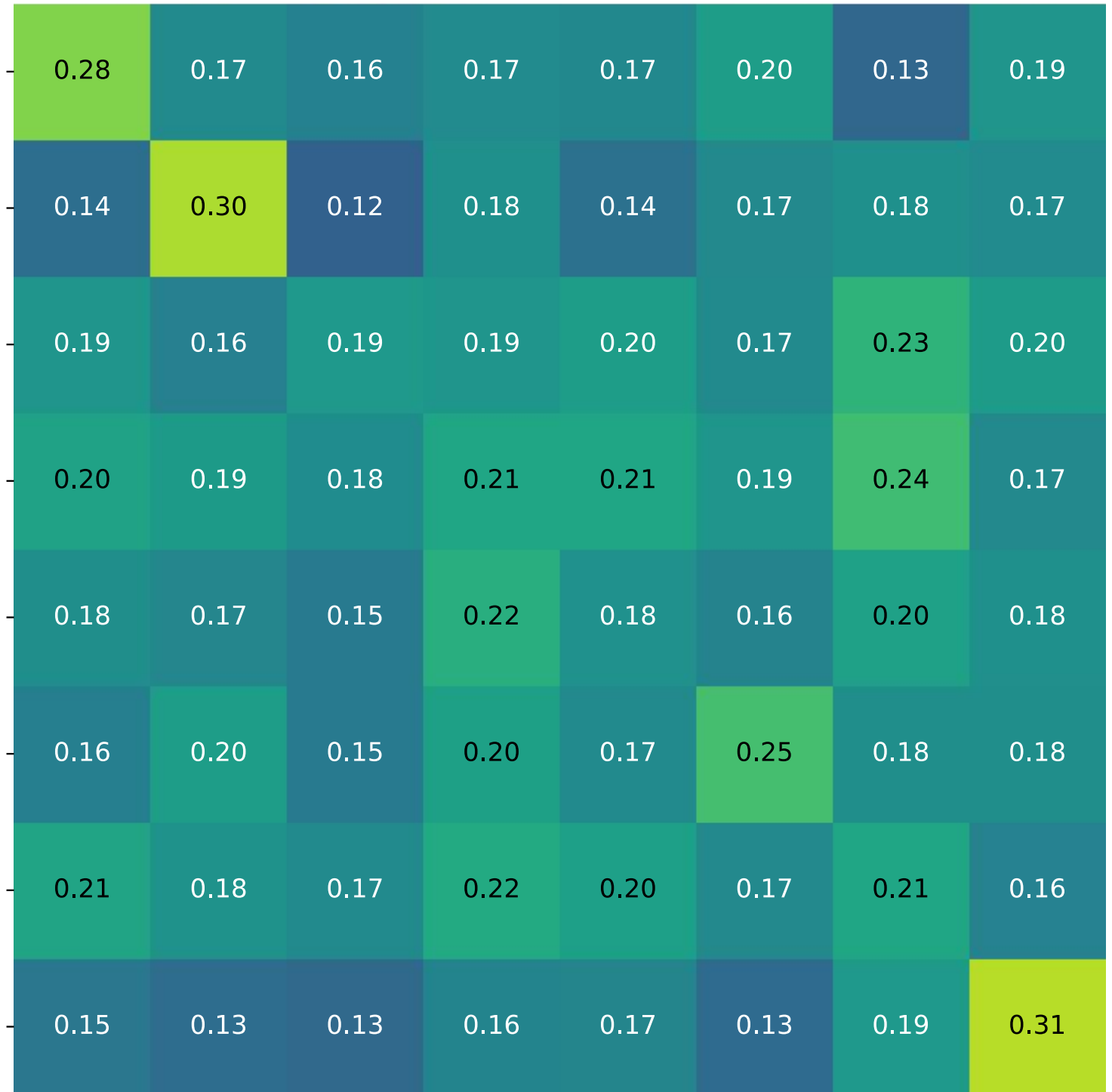
(Ukrainian) Рожевий квіт гілки фруктового дерева на тлі голубого з хмарами неба

(Persian) می‌دق نامت‌خاس کی ل‌خاد زا ی‌ی‌امن و ه‌ب‌اه ه‌رج‌ن‌پ‌ی‌ال ه‌بال زا ه‌ک‌رون و ه‌ب‌ور‌خم و تس‌ا ه‌دی‌بات ل‌خاد

(Bengali) চারদিকে সবুজ গাছপালা আছে এবং মাঝে মুরগী দেখা যাচ্ছে

(Greek) Φόρμα πληροφοριών

(Indonesian) Salah satu dim sum yang dinamakan xiaolongbao disajikan dengan beragam varian warna



# OpenCLIP - Multilingual

Median of matches (diagonal): 0.235

Median of misses (off-diagonal): 0.045

Median Match-to-Miss Delta: 0.190



(English) A young man smiling and posing.

(French) 7 personnes sur une boué jaune et orange faisant du rafting en vertical dans l'eau

(Korean) 타투가 들어있는 사탕 케이스 껍질

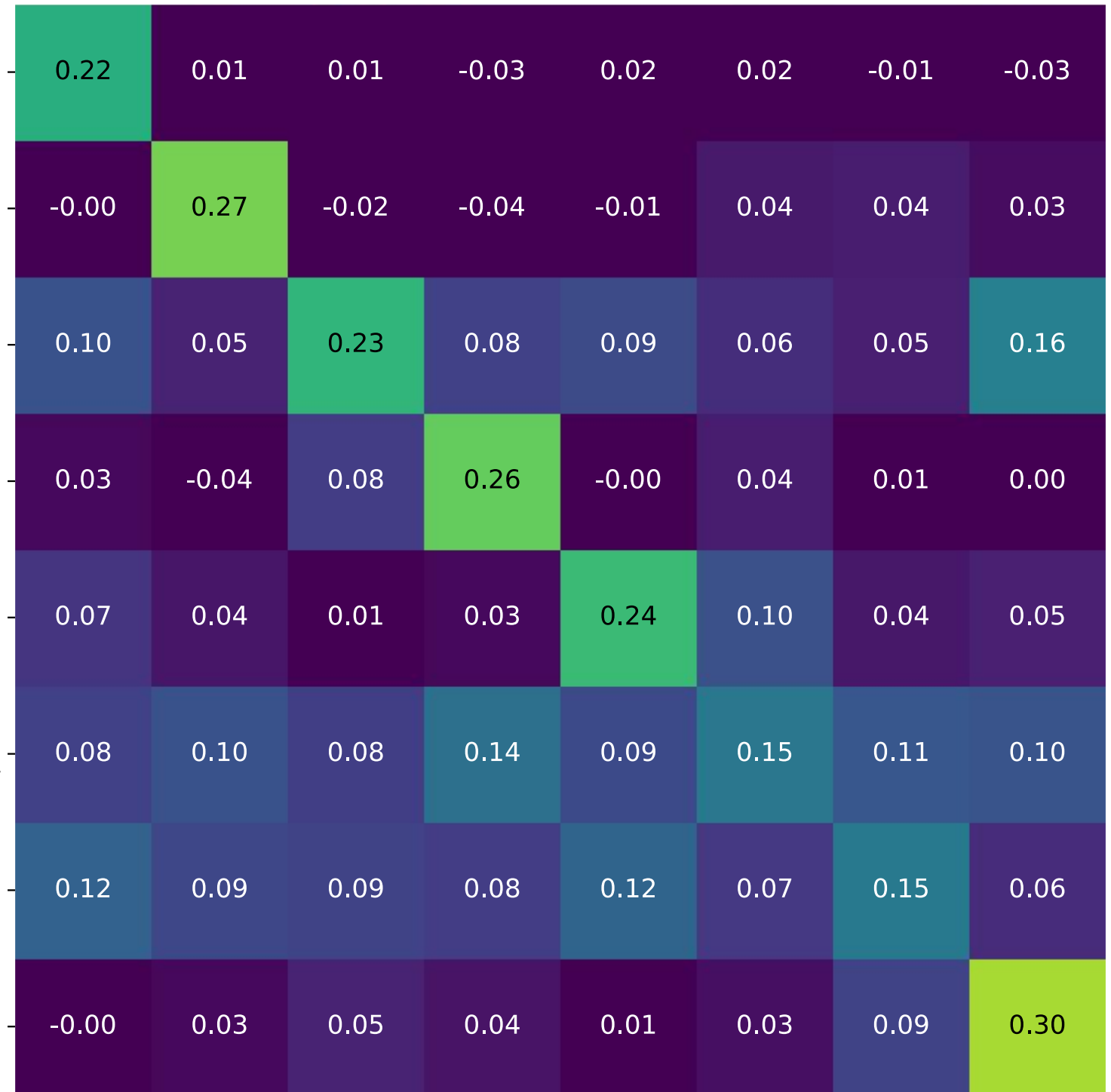
(Ukrainian) Рожевий квіт гілки фруктового дерева на тлі голубого з хмарами неба

(Persian) میمى دق نامتخاس كى لخدازا یىامن هب اه هرجنپ یال هبال زا هكرون و هبورخم و تس ا هدیبات لخدازا

(Bengali) চারদিকে সবুজ গাছপালা আছে এবং মাঝে মুরগী দেখা যাচ্ছে

(Greek) Φόρμα πληροφοριών

(Indonesian) Salah satu dim sum yang dinamakan xiaolongbao disajikan dengan beragam varian warna



# Results from CLIP, Search term: “거실에서 이야기를 나누는 두 사람.” “Two people talking in a living room.” in Korean

Score: 0.2717



Score: 0.2698



Score: 0.2637



Score: 0.2598



Score: 0.2593



# Results from OpenCLIP, Search term: “거실에서 이야기를 나누는 두 사람.” “Two people talking in a living room.” in Korean

Score: 0.2771



Score: 0.2100



Score: 0.2067



Score: 0.1974

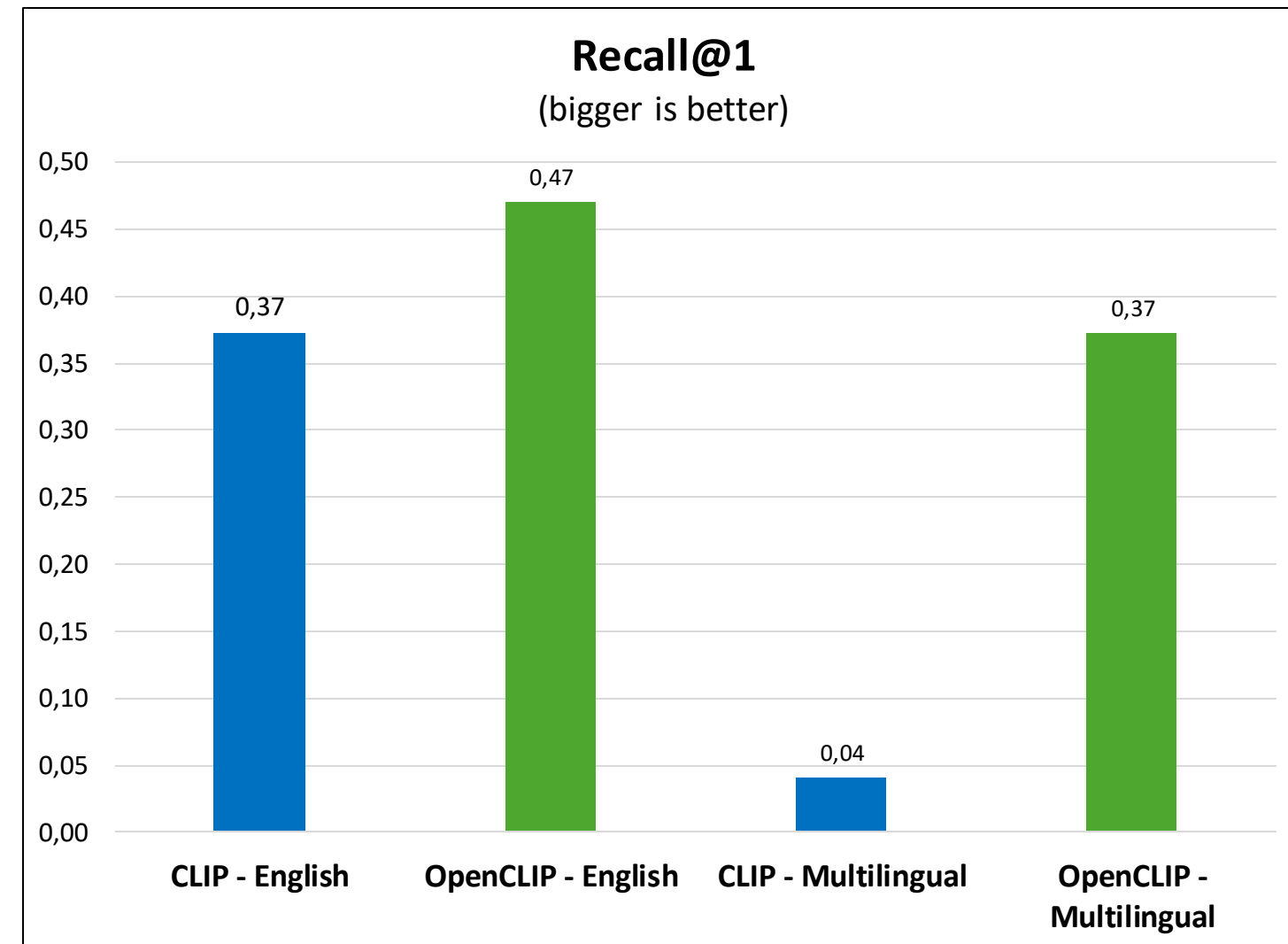
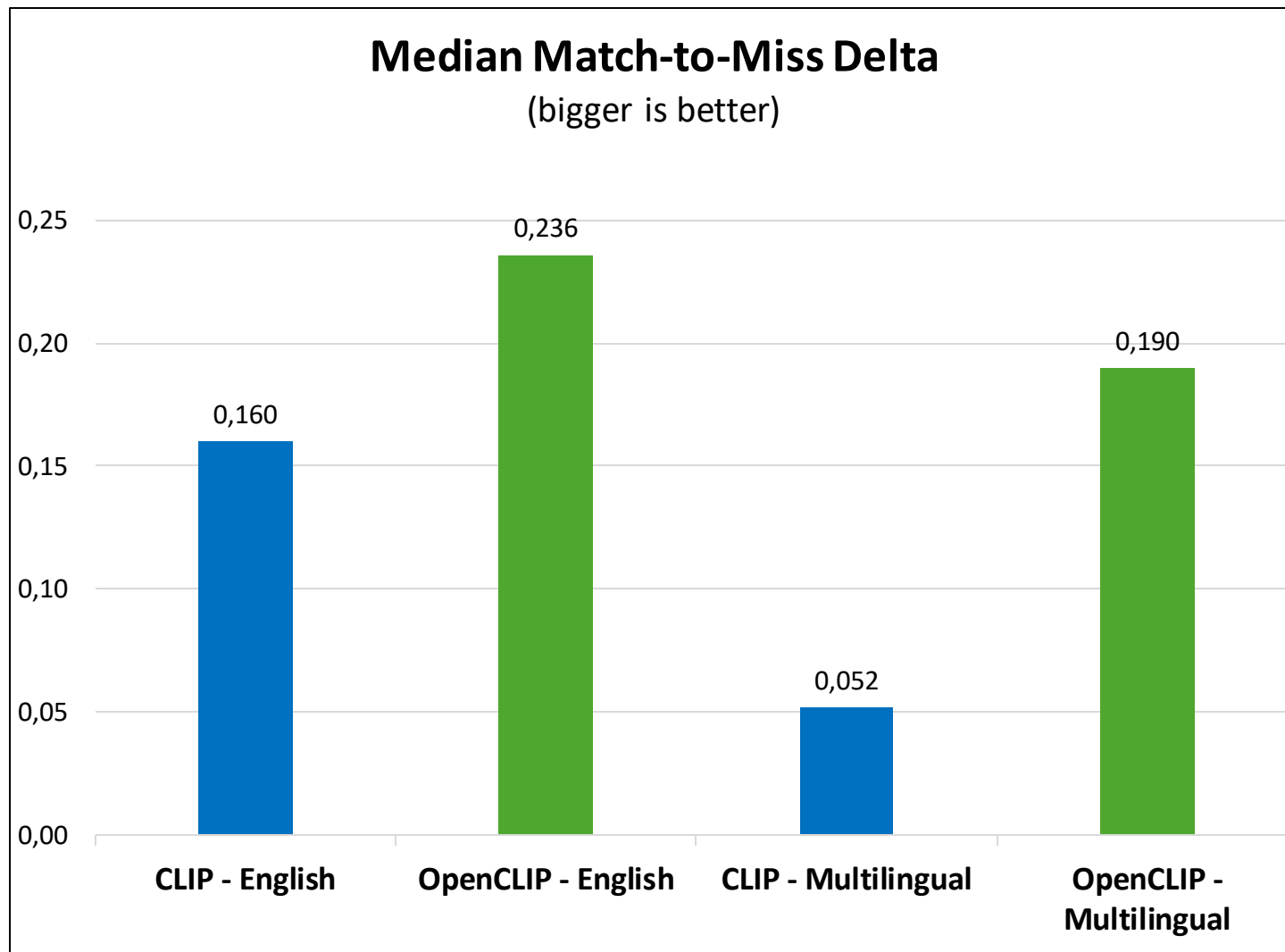


Score: 0.1962





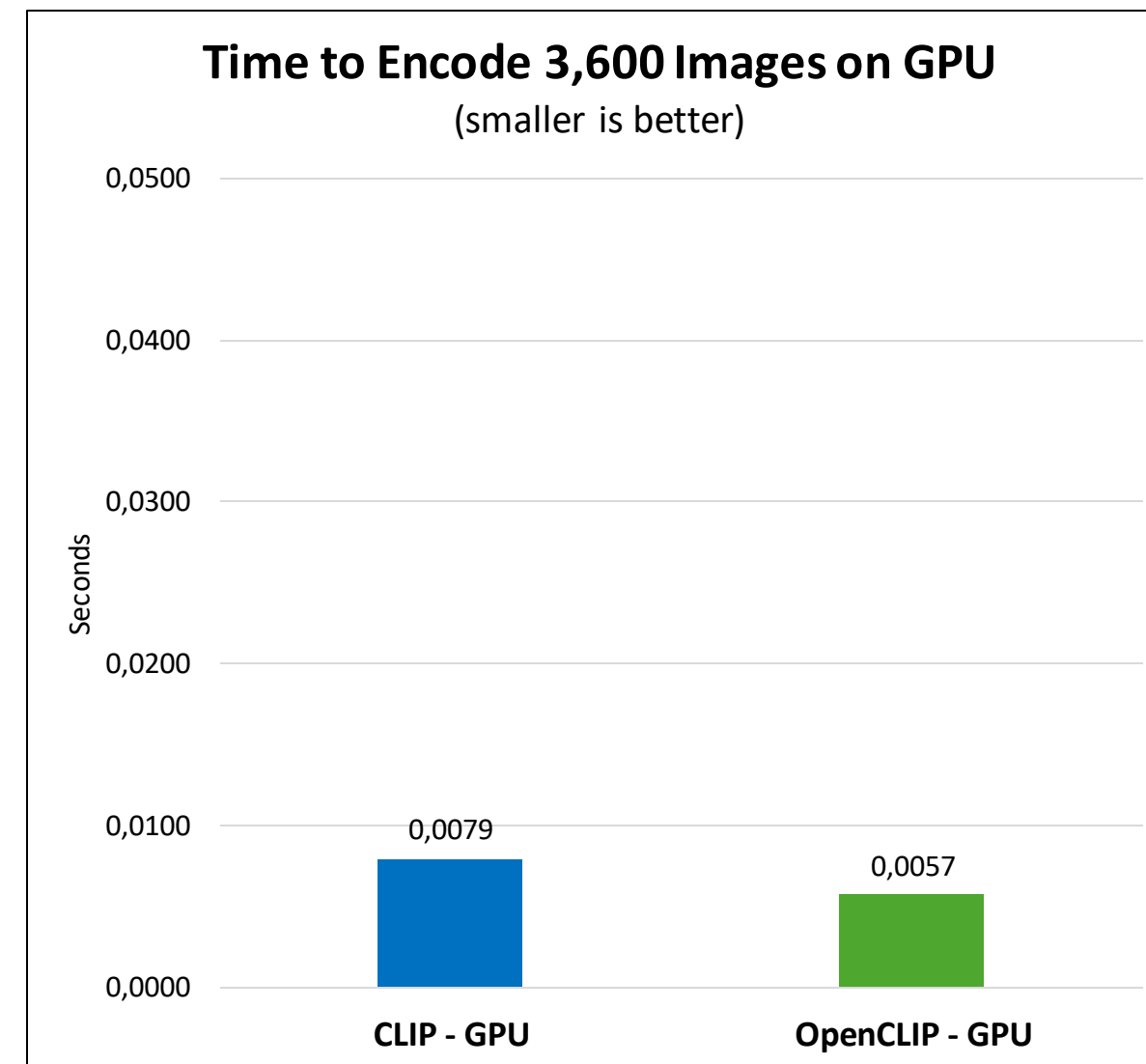
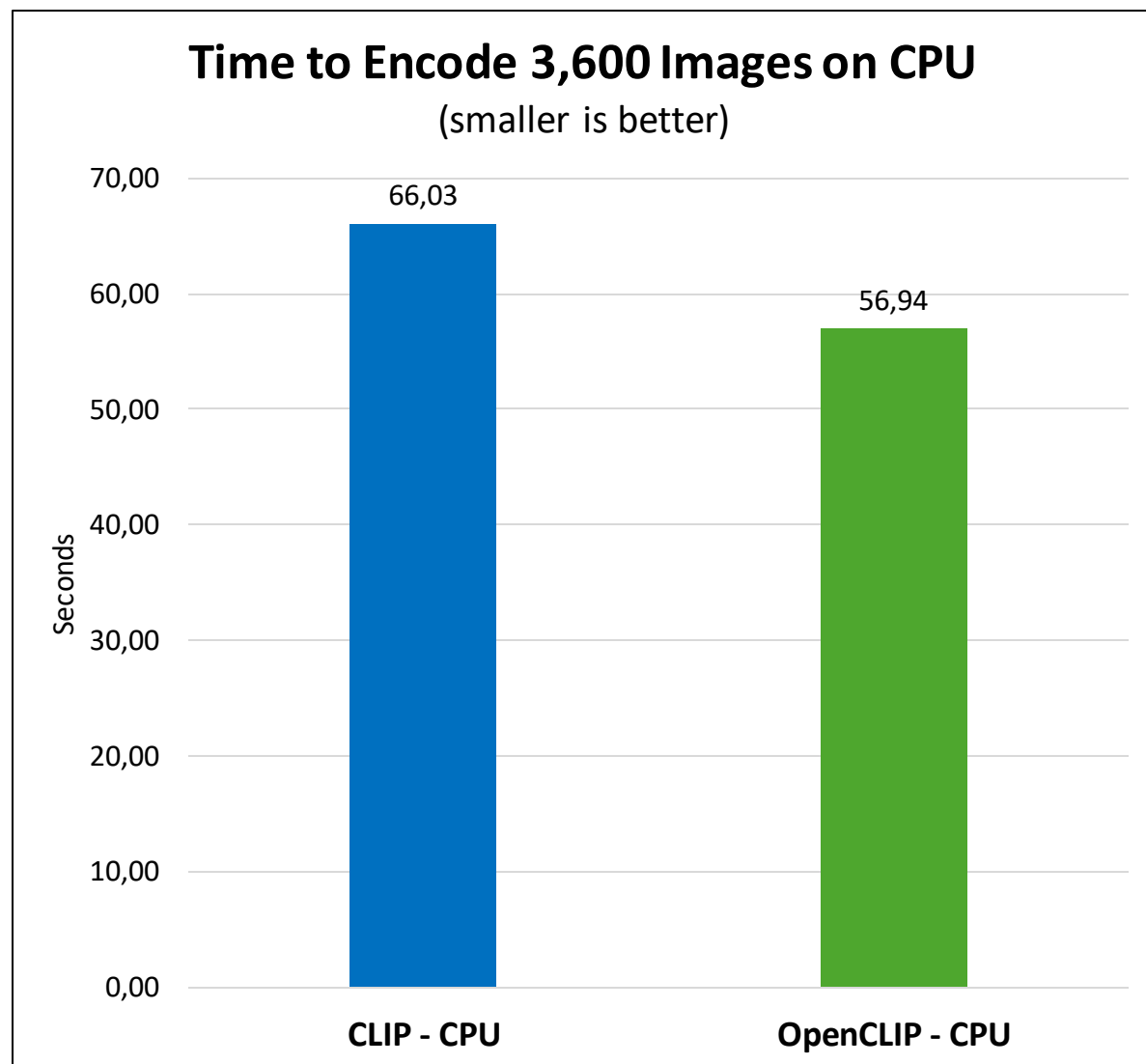
# Accuracy Results of CLIP and OpenCLIP



Dataset: Google's Crossmodal-3600  
Recall@1 results from LAION



# Performance Results of CLIP and OpenCLIP



Dataset: Google's Crossmodal-3600

CPU System: Lenovo ThinkPad P16, Core i9, 24 Cores, 2.3 GHz, 64 GB RAM

GPU System: NVIDIA RTX A5500, 16 GB RAM



# CHALLENGES AND NEXT STEPS



## Challenges:

- Scalability and performance optimization for large volumes of media data
- Ensuring high accuracy across languages
- Extending search to other media types like video and audio

## Next Steps:

- Optimizing encoding speeds with model quantization
- Further fine-tuning models for specific use cases



# CONCLUSION

- Semantic Content Discovery is a method of searching for and retrieving media files
- A database can be used to store the resulting embeddings
- Search using unstructured text queries
- LAION's OpenCLIP outperforms OpenAI's CLIP on many criteria
- OpenCLIP works well with multiple languages





QR code for Avid Ada



robgon-art



@Rob\_Gonsalves\_



medium.com/@robgon