

# Week 4 Report

Thomas Deatherage  
NFHM/BioCosmos  
13 September 2024

## Time Slot

### 1) What progress did you make in the last week?

- Met with the GaTech Team (Vy, Roman and myself on Wednesday evening)
- Met with GaTech + UoF collaborators.
- Continued last week's WIP (DB versioning + embedding tags for segmenting experiments)

### 2) What are you planning on working on next?

- Per today's meeting, I'm going to focus on Florence-2's trait grounding and trait referring capabilities.

### 3) Is anything blocking you from getting work done?

- I'll be traveling for a wedding the latter half of the week, so I may not have a report for next week.

## Abstracts & Summaries

### 1) Re-TASK: Revisiting LLM Tasks from Capability, Skill, and Knowledge Perspectives

*Conference/Journal:* Not stated

*Abstract:* As large language models (LLMs) continue to scale, their enhanced performance often proves insufficient for solving domain-specific tasks. Systematically analyzing their failures and effectively enhancing their performance remain significant challenges. This paper introduces the Re-TASK framework, a novel theoretical model that Revisits LLM Tasks from cApability, Skill, Knowledge perspectives, guided by the principles of Bloom's Taxonomy and Knowledge Space Theory. The Re-TASK framework provides a systematic methodology to deepen our understanding, evaluation, and enhancement of LLMs for domain-specific tasks. It explores the interplay among an LLM's capabilities, the knowledge it processes, and the skills it applies, elucidating how these elements are interconnected and impact task performance. Our application of the Re-TASK framework reveals that many failures in domain-specific tasks can be attributed to insufficient knowledge or inadequate skill adaptation. With this insight, we propose structured strategies for enhancing LLMs through targeted knowledge injection and skill adaptation.

Specifically, we identify key capability items associated with tasks and employ a deliberately designed prompting strategy to enhance task performance, thereby reducing the need for extensive fine-tuning. Alternatively, we fine-tune the LLM using capability-specific instructions, further validating the efficacy of our framework. Experimental results confirm the framework's effectiveness, demonstrating substantial improvements in both the performance and applicability of LLMs.

**Summary:** This paper presents the Re-TASK framework, which helps improve large language models (LLMs) for specific tasks by analyzing their capabilities, knowledge, and skills. It identifies reasons for their failures and suggests targeted strategies to enhance their performance without needing extensive adjustments.

**Link:** <https://arxiv.org/abs/2408.06904>

**Code:** N/A

**Relevance:** Part of our long term vision is to have an LLM interface capable of delegating queries to appropriate, well-defined tasks. This paper is part of my research into that.

## Scripts and Code Blocks

Code contributions can be found at this PR: <https://github.com/BioCosmos-AI/BioCosmos/pull/2>

This week is a continuation of last week. Primary accomplishments:

- 1) Finally got DB versioning working with Flyway. I'm going to have Roman or Vy test this on their machines before I merge.
- 2) Support end-to-end arbitrary tagging from the generate-embedding stage to API search

Some highlights:

```

1 + #!/bin/bash
2 + set -e
3 +
4 + # Function to check if PostgreSQL is ready
5 + pg_isready() {
6 +     PGPASSWORD=$POSTGRES_PASSWORD psql -h localhost -U "$POSTGRES_USER" -d
7 +     "$POSTGRES_DB" -c "SELECT 1" > /dev/null 2>&1
8 + }
9 +
10 + # Start PostgreSQL
11 + /usr/local/bin/docker-entrypoint.sh postgres &
12 +
13 + # Wait for PostgreSQL to be ready
14 + until pg_isready; do
15 +     echo "Waiting for PostgreSQL to be ready..."
16 +     sleep 2
17 + done
18 +
19 + echo "PostgreSQL is ready. Creating extensions..."
20 + # Create extensions
21 + psql -v ON_ERROR_STOP=1 --username "$POSTGRES_USER" --dbname "$POSTGRES_DB" <<-
22 +     EOSQL
23 +     CREATE EXTENSION IF NOT EXISTS postgis;
24 +     CREATE EXTENSION IF NOT EXISTS vector;
25 +     EOSQL
26 +
27 + echo "Extensions created. Running Flyway migrations..."
28 +
29 + # Run Flyway migrations with baseline
30 + flyway -url=jdbc:postgresql://localhost:5432/$POSTGRES_DB -user=$POSTGRES_USER -
31 +     password=$POSTGRES_PASSWORD -baselineOnMigrate=true -
32 +     locations=filesystem:/flyway/sql migrate
33 +
34 + echo "Flyway migrations completed."
35 +
36 + # Keep the container running
37 + wait

```

New entrypoint into Postgres

```

@router.post("/search", response_model=SearchResponse)
v async def search(
    request: Request,
    search_param: Optional[str] = Form(None),
    image: Optional[UploadFile] = File(None),
    limit: int = Query(30, ge=1, le=100),
    embed_version: str = Query("default", max_length=512),
    session: AsyncSession = Depends(get_session)
) -> SearchResponse:
    try:
        if search_param is None and image is None:
            raise HTTPException(status_code=400, detail="At least one of search_param or image must be provided.")

        search_vector = await process_input(request, search_param, image)

        query = select(SearchRecord).where(
            SearchRecord.embed_version == embed_version
        ).order_by(
            func.l2_distance(
                SearchRecord.embedding,
                cast(search_vector, Vector)
            )
        ).limit(limit)

        results = await session.execute(query)
        records = results.scalars().all()

        payload = [create_record_payload(record) for record in records if record.external_media_uri is not None]

        return SearchResponse(
            search_param=search_param,
            filename=image.filename if image else None,
            record_count=len(payload),
            records=payload,
            embed_version=embed_version
        )

```

The /search endpoint now accepts a query parameter “embed\_version” to filter tagged rows.

## Flow Charts/Diagrams

Nothing new flow-chart wise to show.

## Documentation

I added documentation to the repo’s README

## Database Versioning.

We use Flyway to version the Postgres DB. To make a new migration/version, add your migration changes to a file named per the [Flyway Naming Conventions](#) to the `postgres/migrations` directory. Then rebuild the docker container.

## Running and tagging experiments

When generating embeddings, you have the option to tag rows in Postgres with an arbitrary string of text. This is useful for running experiments and comparing results. The string will be stored in the Postgres database and can be used to identify the embedding later.

### Example

```
$ bin/ingest_embed --embed_version "thomas_experiment_22" --pretrained "ViT-B-32" --model "laion2b_s34b_b79k"
```

- From the workbench of Redis Insight, pass a simple search string to the `embedder` queue:
  - `LPUSH embedder "{}"`

From the API, you can limit results to a specific "embed\_version" tag:

```
curl -X 'POST' \
  'http://localhost:8080/api/search?limit=30&embed_version=thomas_experiment_22' \
  -H 'accept: application/json' \
  -H 'Content-Type: multipart/form-data' \
  -F 'search_param=butterfly'

{
  "search_param": "butterfly",
  "filename": null,
  "record_count": 100,
  "records": [ . . . ],
  "embed_version": "thomas_experiment_22"
}
```

Note this only works for Open CLIP generated embeddings at the moment. (TODO: Update search query to select correct embedding model by "embed\_version").

## Results Visualization + Proof of Work

When generating embeddings, we can now tag an arbitrary "embed\_model".

```
$ bin/ingest_embed --embed_model=thomas_experiment_1
```

And on the read-side of the equation, we can filter our focus to a specific tagged version of generated embeddings:

Curl

```
curl -X 'POST' \
'http://localhost:8080/api/search?limit=30&embed_version=thomas_experiment_1' \
-H 'accept: application/json' \
-H 'Content-Type: multipart/form-data' \
-F 'search_param=string'
```

Request URL

```
http://localhost:8080/api/search?limit=30&embed_version=thomas_experiment_1
```

Server response

Code	Details
200	<p>Response body</p> <pre>{   "search_param": "string",   "filename": null,   "record_count": 30,   "records": [     {       "id": 728,       "scientific_name": "Iophanus pyrrihas",       "common_name": "blues and coppers; butterflies; butterflies and moths; winged insects; insects; hexapods; arthropods; animals",       "name": "blues and coppers; butterflies; butterflies and moths; winged insects; insects; hexapods; arthropods; animals",       "description": "blues and coppers; butterflies; butterflies and moths; winged insects; insects; hexapods; arthropods; animals",       "image_source_name": "",       "specimen_source_name": "",       "external_id": "994431c8-de4e-42de-b290-723dcae03c41",       "media_url": "https://images.collections.yale.edu/iiif/2/ypm:021958a1-efe1-4059-90d6-864a3235de0d/full/11920,1920/0/default.jpg",       "specimen_id": "142920df-7b0d-4496-91a7-42be0c2580d9",       "recorded_by": "Eduardo C. Welling",       "collection_date": "1966-11-04",       "source": "iDigBio",       "latitude": 18.534509,       "longitude": -98.943676,       "model": "ViT-B-32",       "pretrained": "laion2b_s34b_b79k",       "embed_version": "thomas_experiment_1"     }   ], }</pre> <p>Response headers</p> <pre>content-length: 26966 content-type: application/json date: Fri, 13 Sep 2024 20:41:54 GMT server: uvicorn</pre>

Responses

Code	Description	Links
------	-------------	-------

## Next Week's Proposal

- Merge this week's work
- Focus on Florence-2 and training groundings/referrings