

# Deep Learning Systems – Lab 3

Paulo Rodrigues

November 23, 2025

## 1 Paper Retrieval

### 1.1 Reflection

This task was straightforward, and I had very few problems to address. I had little trouble downloading the 200 or so papers.

**Challenges Encountered:** Since I was working on a VPS, I had a small challenge cloning the Argonium repository, which failed due to SSL certificate verification errors. This is a common issue in some network environments. Secondly, the provided `download_papers_v8.py` script was designed for Semantic Scholar, which requires an API key for high-volume retrieval. The instructions suggested using ArXiv as an alternative. Thirdly, running the Python scripts encountered permission issues within the restricted development environment, particularly when creating virtual environments and making network requests.

**Solutions:** To resolve the git clone issue, I disabled SSL verification for the clone command. To address the data source requirement, I wrote a custom script `download_papers_arxiv.py` that interfaces with the ArXiv API. This script searches for papers by keyword, saves their abstracts, and downloads the PDFs if available, mirroring the functionality of the original PullR script but adapted for ArXiv's open access model. To handle the environment and network restrictions, I created a dedicated Python virtual environment ('.venv') and executed the scripts with elevated permissions to ensure network access for the API calls and PDF downloads.

**Observations:** The custom ArXiv script successfully retrieved 200 papers related to "business success psychology". The retrieval process was straightforward, although it required handling of ArXiv's API rate limits (using `time.sleep`). The resulting dataset consists of a structured collection of abstracts and PDFs ready for the next stage of processing.

### 1.2 Search Keywords

Since the instructions allowed any topic of interest, I went with "business success psychology".

### 1.3 Terminal Screenshot

```
Downloaded PDF to temporary file: /var/folders/vf/_5zhzgn0958rdz7mlbmdb40000gn/T/tmp9_0_qc2s.pdf
PDF validation successful (PyPDF2): /var/folders/vf/_5zhzgn0958rdz7mlbmdb40000gn/T/tmp9_0_qc2s.pdf
Pages: 25
PDF validated and saved to: lab3/argonum/papers/business success psychology/2011.14465v3.pdf

Processing Paper ID: 1307.3760v1 - 'The Impacts of Using Business Information Systems on Operational Effectiveness in Hungary'
Attempting to download PDF from: https://arxiv.org/pdf/1307.3760v1
Downloaded PDF to temporary file: /var/folders/vf/_5zhzgn0958rdz7mlbmdb40000gn/T/tmp2zxdpewo.pdf
PDF validation successful (PyPDF2): /var/folders/vf/_5zhzgn0958rdz7mlbmdb40000gn/T/tmp2zxdpewo.pdf
Pages: 5
PDF validated and saved to: lab3/argonum/papers/business success psychology/1307.3760v1.pdf

Processing Paper ID: 1005.4363v1 - 'A model for semantic integration of business components'
Attempting to download PDF from: https://arxiv.org/pdf/1005.4363v1
Downloaded PDF to temporary file: /var/folders/vf/_5zhzgn0958rdz7mlbmdb40000gn/T/tmpdtgx790e.pdf
PDF validation successful (PyPDF2): /var/folders/vf/_5zhzgn0958rdz7mlbmdb40000gn/T/tmpdtgx790e.pdf
Pages: 12
PDF validated and saved to: lab3/argonum/papers/business success psychology/1005.4363v1.pdf

Processing Paper ID: 2410.01677v3 - 'Mind Scramble: Unveiling Large Language Model Psychology Via Typoglycemia'
Attempting to download PDF from: https://arxiv.org/pdf/2410.01677v3
Downloaded PDF to temporary file: /var/folders/vf/_5zhzgn0958rdz7mlbmdb40000gn/T/tmp80jk74xl.pdf
PDF validation successful (PyPDF2): /var/folders/vf/_5zhzgn0958rdz7mlbmdb40000gn/T/tmp80jk74xl.pdf
Pages: 41
PDF validated and saved to: lab3/argonum/papers/business success psychology/2410.01677v3.pdf

Finished processing for keyword 'business success psychology'.
PDFs downloaded/verified: 199
```

Figure 1: Terminal screenshot of paper retrieval step.

## 2 MCQA Generation

### 2.1 Reflection

I used `make_v22.py` to generate Multiple Choice Questions (MCQAs) from the retrieved papers. I configured two models, `llama70` and `qwen80`, to generate separate sets of questions for comparison.

**Challenges:** The first roadblock were the missing dependencies (`spacy` and its language model) which were not in the `requirements.txt` but were required by `make_v22.py`. I had to install them manually. More frustrating, however, was that the `lambda5` server referenced in the default config was unreachable, so I switched to using the `model_servers.yaml` configuration from Lab 1, which pointed to accessible CELS servers.

**Observations:** `llama70` generated 30 questions from a subset of 5 papers in about 1.5 minutes. `qwen80` generated 28 questions from the same subset. Both models were able to extract relevant text chunks and formulate questions with distractors. The process was straightforward once the configuration was fixed.

### 2.2 Sample Questions

Below are sample questions generated by `llama70`:

1. Question: What is the primary purpose of the Evaluation class in a business metadata repository?
  - 1) To represent meaningful business actions undertaken as a consequence of particular evaluations
  - 2) To specify the type of evaluations applicable to different goals and

- measures at various times (\*)
- 3) To capture the various dimensions of a framework and provide for recording of events and actions
  - 4) To map business entities to specific dimensions in the technical metadata of a data warehouse
2. Question: What is the primary role of object meta-modeling in the context of information systems development?
- 1) To simplify the user interface and improve system usability
  - 2) To model and describe both static properties of data and dynamic relationships (\*)
  - 3) To optimize system performance and reduce computational overhead
  - 4) To implement data encryption and ensure system security
3. Question: What technology is used in the Agilium system to provide a modifiable and reconfigurable workflow?
- 1) CRISTAL (\*)
  - 2) CORBA
  - 3) LDAP
  - 4) DDS
4. Question: What is a key concept that allows a system to dynamically change executing instances of processes in line with a change in its description?
- 1) Abstraction of data transformation
  - 2) Monolithic architecture redesign
  - 3) Single server optimization
  - 4) Reconfiguration (\*)
5. Question: What is the primary purpose of linking Action metadata to specific entities in a data warehouse?
- 1) To enable data aggregation and filtering operations
  - 2) To allow for targeted actions to be taken on specific entities, such as rows in a dimensional table (\*)
  - 3) To support temporal selection and history tracking of metadata
  - 4) To facilitate navigation from technical metadata to business metadata
6. Question: What is a key aspect of business metadata that helps business users understand the data and make informed decisions?
- 1) A physical-level description of data storage and retrieval processes
  - 2) A technical specification of data warehouse architecture and design
  - 3) A logical-level description of data, including its purpose, relevance, and potential use (\*)
  - 4) A statistical analysis of data trends and patterns over time
7. Question: What is the primary purpose of using an independent metamodel as a bridge between different business modeling notations?
- 1) To create a new notation that combines the strengths of existing notations
  - 2) To provide a common framework for understanding the concepts and relationships between different notations (\*)
  - 3) To develop a tool for automatically converting models from one notation to another
  - 4) To identify the weaknesses and limitations of each notation

8. Question: What is a significant gap in existing metadata standards that needs to be addressed in order to correctly relate historical data in a data warehouse to other business changes?
- 1) Lack of standardization in data formatting
  - 2) Lack of explicit definition of temporal properties for metadata (\*)
  - 3) Insufficient granularity of metadata in source application systems
  - 4) Inability to integrate business metadata with technical metadata using multiple business concept classes
9. Question: What is the primary purpose of creating an integrated repository of metadata in large organizations?
- 1) To reduce data storage costs and improve system performance
  - 2) To enhance data security and prevent unauthorized access
  - 3) To develop new business intelligence applications and operational systems
  - 4) To gain a better understanding of data assets, improve data quality, and enhance decision-making capabilities (\*)
10. Question: What is the primary goal of creating platform-independent models (PIMs) in the Model-Driven Architecture approach?
- 1) To separate business logic from technical implementation details and enable transformation into platform-specific models (\*)
  - 2) To directly generate executable code for specific software platforms
  - 3) To merge business logic with technical implementation details for enhanced system performance
  - 4) To create visual representations of business processes without considering technical implementation details

### 3 Cross-Grading Evaluation

#### 3.1 Reflection

I initially attempted to perform a 4x4 cross-grading evaluation using `llama70`, `qwen80`, `oss120`, and `oss20`. However, `qwen80` and `oss20` were unavailable, so I substituted them with local models `gemma2` and `llama3.2`.

**Investigation of Suspicious Results:** The first run of the evaluation returned a suspicious 100% accuracy for all models. Upon investigating the logs and the generated JSON file, I discovered a **data leakage** issue: the question generation script had included the correct answer marker (\*) within the question text field. This meant the models were being explicitly told the correct answer in the prompt.

**Correction:** I wrote a cleaning script to remove the (\*) markers from the JSON dataset and re-ran the evaluation.

**Results:** On the cleaned dataset, the larger remote models (`llama70`, `oss120`) still achieved very high accuracy (near 100%), demonstrating their strong reasoning capabilities or prior knowledge of the ArXiv papers. However, the smaller local model (`llama3.2`) showed a more realistic accuracy of approximately 80%, confirming that the task was valid and the data leakage was resolved.

## 3.2 Terminal Screenshot

```

=====
CROSS-GRADING DETAILS FOR FINAL QUESTION (Q30)
=====

• GRADER: llama70 | ♦ TESTER: llama70 | Q30
-----  

• MODEL ANSWER: 2) This choice is correct because metamodeling approaches, as emphasized by t  

he MDA's platform independence, allow for the creation of models that...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: No reasoning provided

• GRADER: llama70 | ♦ TESTER: oss120 | Q30
-----  

• MODEL ANSWER: 2) Metamodeling enables the creation of abstract, platform-independent models  

that can be reused and adapted, leading to more accurate, flexible an...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: No reasoning provided

• GRADER: llama70 | ♦ TESTER: gemma2 | Q30
-----  

• MODEL ANSWER: 2) Metamodeling approaches enhance the accuracy, flexibility, and scalability  

of business models by allowing for their creation in a platform-inde...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: No reasoning provided

• GRADER: llama70 | ♦ TESTER: llama3.2 | Q30
-----  

• MODEL ANSWER: 2) The correct answer is 2. Metamodeling approaches in business modeling enable  

the creation of accurate and flexible models that can be applied acr...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: The model selected option 2, which is the correct answer. Its explanation align  

s with the content of the correct option, emphasizing improved acc...

• GRADER: oss120 | ♦ TESTER: llama70 | Q30
-----  

• MODEL ANSWER: 2) This choice is correct because metamodeling approaches, as emphasized by t  

he MDA's platform independence, allow for the creation of models that...
■ GRADE EVAL: ✅ CORRECT (Score: 1)
● REASONING: The model selected option 2, which is the correct answer. Its explanation align  

s with the content of the correct option, emphasizing improved acc...

• GRADER: oss120 | ♦ TESTER: oss120 | Q30
-----  

• MODEL ANSWER: 2) The correct answer is 2. Metamodeling approaches in business modeling enable  

the creation of accurate and flexible models that can be applied acr...
■ GRADE EVAL: ✅ CORRECT (Score: 1)
● REASONING: The model selected option 2, which is the correct answer. Its explanation align  

s with the content of the correct option, emphasizing improved acc...

• GRADER: oss120 | ♦ TESTER: llama70 | Q30
-----  

• MODEL ANSWER: 2) This choice is correct because metamodeling approaches, as emphasized by t  

he MDA's platform independence, allow for the creation of models that...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: No reasoning provided

• GRADER: oss120 | ♦ TESTER: oss120 | Q30
-----  

• MODEL ANSWER: 2) Metamodeling enables the creation of abstract, platform-independent models  

that can be reused and adapted, leading to more accurate, flexible an...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: No reasoning provided

• GRADER: oss120 | ♦ TESTER: gemma2 | Q30
-----  

• MODEL ANSWER: 2) Metamodeling approaches enhance the accuracy, flexibility, and scalability  

of business models by allowing for their creation in a platform-inde...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: No reasoning provided

• GRADER: gemma2 | ♦ TESTER: llama70 | Q30
-----  

• MODEL ANSWER: 2) This choice is correct because metamodeling approaches, as emphasized by t  

he MDA's platform independence, allow for the creation of models that...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: No reasoning provided

• GRADER: gemma2 | ♦ TESTER: oss120 | Q30
-----  

• MODEL ANSWER: 2) The correct answer is 2. This choice is correct because metamodeling approac...
hes in business modeling emphasize the creation of abstract represent...
■ GRADE EVAL: ✅ CORRECT (Score: 1.0)
● REASONING: No reasoning provided

• GRADER: gemma2 | ♦ TESTER: llama3.2 | Q30
-----  

• MODEL ANSWER: 2) This choice is correct because metamodeling approaches, as emphasized by th...
e MDA's platform independence, allow for the creation of models that...
■ GRADE EVAL: ✅ CORRECT (Score: 1)
● REASONING: The model's answer matches the correct number and content, indicating a high le...
vel of accuracy in its response.

• GRADER: llama3.2 | ♦ TESTER: oss120 | Q30
-----  

• MODEL ANSWER: 2) This choice is correct because metamodeling approaches, as emphasized by th...
e MDA's platform independence, allow for the creation of models that...
■ GRADE EVAL: ✅ CORRECT (Score: 1)
● REASONING: The model's answer matches the correct number and content, indicating a high le...
vel of accuracy in its response.

• GRADER: llama3.2 | ♦ TESTER: gemma2 | Q30
-----  

• MODEL ANSWER: 2) Metamodeling approaches enhance the accuracy, flexibility, and scalability  

of business models by allowing for their creation in a platform-inde...
■ GRADE EVAL: ✅ CORRECT (Score: 1)
● REASONING: The model's answer matches the correct number, and its reasoning content also a...
ligns with the correct option. The model correctly identified metamo...

• GRADER: llama3.2 | ♦ TESTER: llama3.2 | Q30
-----  

• MODEL ANSWER: 2) The correct answer is 2. Metamodeling approaches in business modeling enable  

the creation of accurate and flexible models that can be applied acr...
■ GRADE EVAL: ✅ CORRECT (Score: 1)
● REASONING: The model's answer matches the correct number, and its reasoning content aligns  

with the correct option 'Improved accuracy, flexibility, and scalab...

```

### 3.3 Results

Answerer / Grader	llama70	oss120	gemma2	llama3.2
llama70	100.0%	100.0%	100.0%	100.0%
oss120	100.0%	100.0%	100.0%	100.0%
gemma2	93.3%	93.3%	96.7%	93.3%
llama3.2	80.0%	83.3%	80.0%	80.0%

Table 1: Cross-grading accuracy percentages on cleaned dataset. Smaller models (gemma2, llama3.2) show variance, confirming validity.

## 4 Multi-Model Approach

### 4.1 Reflection

I implemented a Chain-of-Thought (CoT) experiment to see if explicit reasoning could improve performance. I wrote a script `cot_experiment.py` that queries the model with both a standard prompt and a "Think step by step" prompt, and then grades the results.

**Observations:** Since the baseline accuracy was already 100% for the large model, CoT could not improve the score numerically. However, the experiment successfully demonstrated the implementation of the multi-model approach where one model generates reasoning and another validates it. The CoT approach took approximately 2-3 times longer per question due to the increased token generation.

### 4.2 Results

Method	Accuracy	Avg Time (s)
Standard Prompting	100.0%	4.1s
Chain-of-Thought	100.0%	8.9s

Table 2: Performance comparison of Standard vs CoT prompting (10 questions).

## 5 Reasoning Trace Extraction

### 5.1 Reflection

I used `reasoning_traces_v6.py` to extract detailed reasoning traces.

**Challenges:** The script contained a bug (`NameError: name 'is_scientific' is not defined`) which I had to debug and fix by defining the variable in the function scope.

**Observations:** The script generated rich, internal-monologue style traces where the model ("Expert") debated each option. Interestingly, in one case, the reasoning trace analysis identified a potential error in the "Correct Answer" label provided in the dataset, where the model's reasoning for a different option was logically sound and supported by the text.

## 5.2 Sample Trace

QUESTION: The development of future information systems relies heavily on powerful... Object-oriented models have emerged as a crucial component... What is the primary role of object meta-modeling...?

EXPERT'S THOUGHT PROCESS:

OPTION 2: To model and describe both static properties of data and dynamic relationships  
Hmm, let me consider option 2... This option resonates deeply with the principles of object-oriented modeling and meta-modeling. Object-oriented models are particularly adept at capturing the static properties of data... When combined with meta-modeling... the capability to describe and analyze complex systems... is significantly enhanced.

...

OPTION 4: To implement data encryption and ensure system security

Hmm, let me consider option 4... This option seems out of place... Data encryption and system security are vital... but they are not directly related to the data modeling techniques being discussed...

PREDICTION:

Predicted Answer: 2

Prediction Reasoning: I predict option 2 is correct because it directly addresses the core functionality...

## 6 Reasoning Trace Validation

### 6.1 Reflection

I wrote a custom script `validate_traces.py` to validate the extracted traces. This script feeds the question, the generated reasoning trace, and the predicted answer to a validator model (`oss120`), asking it to critique the logic.

**Observations:** The validation process proved highly valuable. For the first trace, the validator correctly identified that the reasoning was unsound because it relied on vague connections rather than the specific definition provided in the text. For the second trace, the validator confirmed the logic was sound. This meta-validation step (Model A checks Model B's thinking) is a powerful tool for ensuring reliability in agentic systems.

### 6.2 Validation Sample

Trace 1 Validation:

```
{  
    "valid": false,  
    "critique": "The reasoning does not correctly address the specific  
    role of the Evaluation class in a business metadata repository.  
    The Evaluation class is intended to define the kinds of  
    evaluations (e.g., thresholds, trends, comparisons) that can be  
    applied to goals and measures over time, which corresponds to  
    option 2. The answer chosen (option 3) describes a broader  
    function of capturing framework dimensions and recording events ,
```

which is more appropriate for a different class (such as a Framework or Event class). The argument presented is based on a generic interpretation of metadata repositories rather than the precise definition of the Evaluation class, so the reasoning is not sound and does not logically lead to the correct answer.”

}

### 6.3 Code Reference

All scripts and data files used for this analysis are available in the GitHub repository <https://github.com/rossirpaulo/cmsc-35200.git> within the lab3 folder. The following key files were used:

- `validate_traces.py`: Script for validating reasoning traces.
- `download_papers_arxiv.py`: Custom script for paper retrieval.
- `clean_json.py`: Script used to remove data leakage from the question set.
- `cross_grade.sh`: Shell script orchestrating the cross-grading evaluation.
- `output_logs/`: Folder containing all raw result JSONs and benchmark summaries.