



# MVP - FOCUS AI

Extension that will help you stay focused!  
Szymon Smagowski, Jerzy Kraszewski

# Table of Contents

- Introduction
- How does it work?
- Data Processing
- Technical Implementation
- Prompt Engineering and results
- Prompt Tuning and results
- LangChain, LangGraph and LangSmith
- Future Steps





# INTRODUCTION

Dataset: (HTML of webpages)

<https://www.kaggle.com/datasets/uciml/identifying-interesting-web-pages/data>



## Content Focus Assistant

A smart application that helps users stay focused on their learning goals:

- Automatically classifies webpage content
- Prevents distractions during work or study
- Analyzes content relevance to user's objectives
- Built using LangChain + OpenAI for intelligent content analysis
- Helps users make better decisions about their time

# How does it work?



## User Browsing

Extension monitors  
webpage content in  
real-time



## AI Analysis

Content is analyzed  
using AI to determine  
relevance



## Decision

System decides to  
allow or block based  
on learning goals

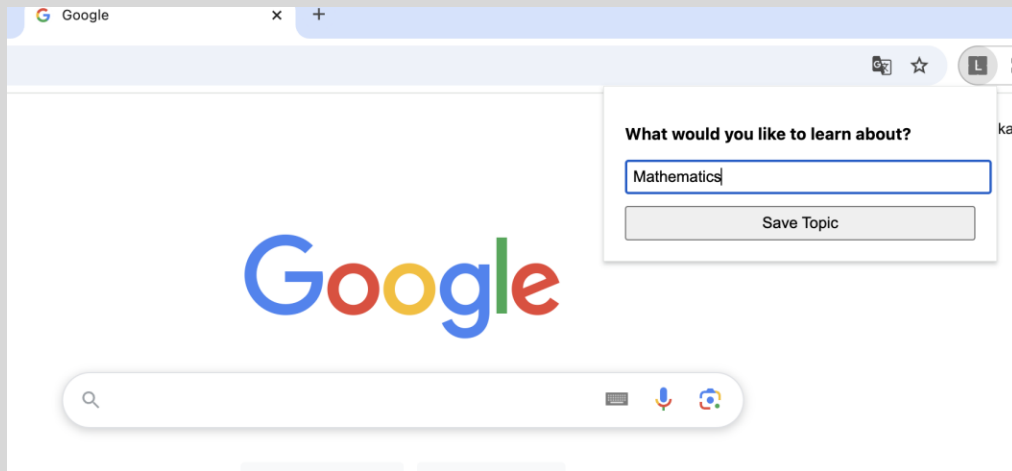


Allow Relevant Content

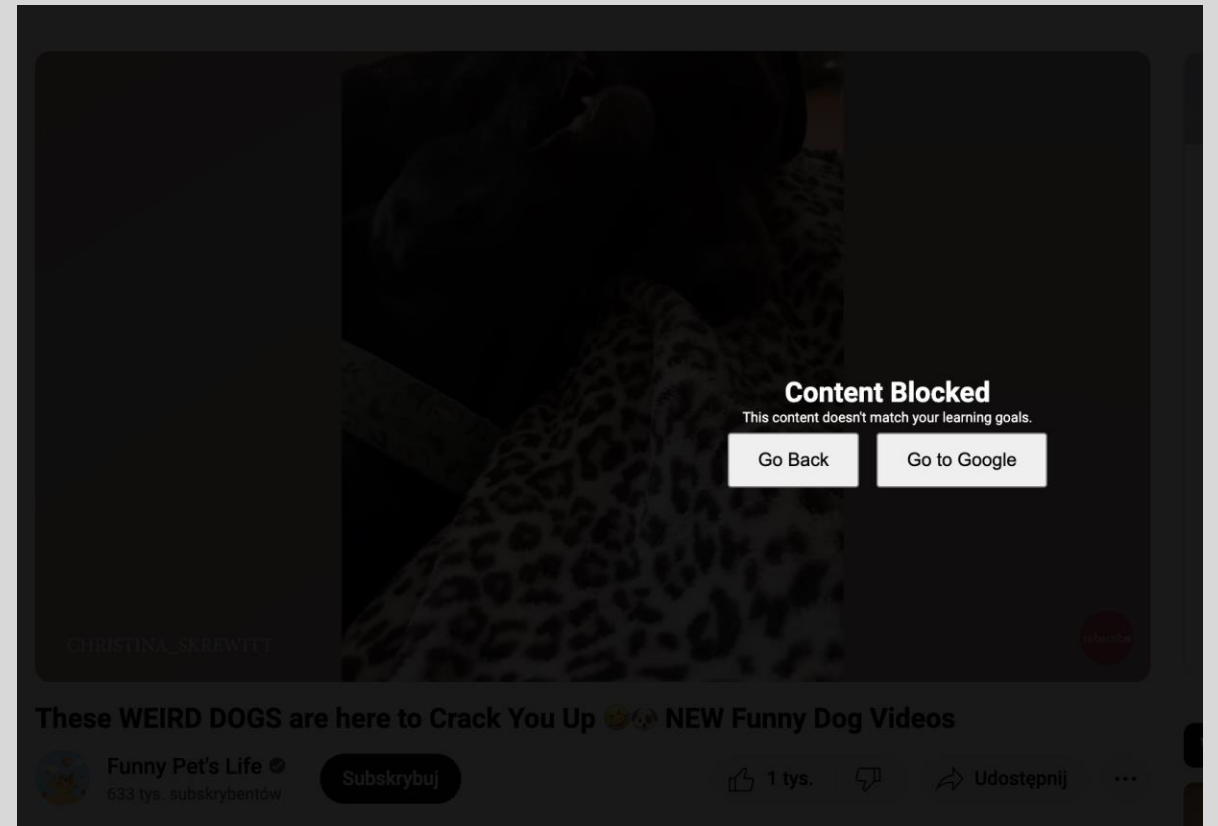


Block Distractions

# How does it work?



Step 1 – user decides on topic of learning



Step 2 – prompted LLM decides if content of the page is relevant for him (in this case it was blocked)



# DATA PROCESSING

## Dataset Overview



**Kaggle Dataset**



### Content Categories

- Medical Articles
- Animal Information
- Scientific Content
- Music
- Others



### Dataset Details

- ~200 HTML Pages
- Various Topics



HTML Files



Markdown Format



# Data Processing 'Critic LLM'

NLP-Focus-AI > processed-data > 📄 ready\_dataset.csv

```
1 markdown_content,matching_label,non_matching_label,matching_status
```

## Dataset Generation



### GPT-4 Label Generation

#### 10 Matching Labels

Labels that describe content

#### 10 Non-Matching Labels

Labels unrelated to content



### Random Sampling

100 randomly selected rows containing:



Markdown Content



Random Label



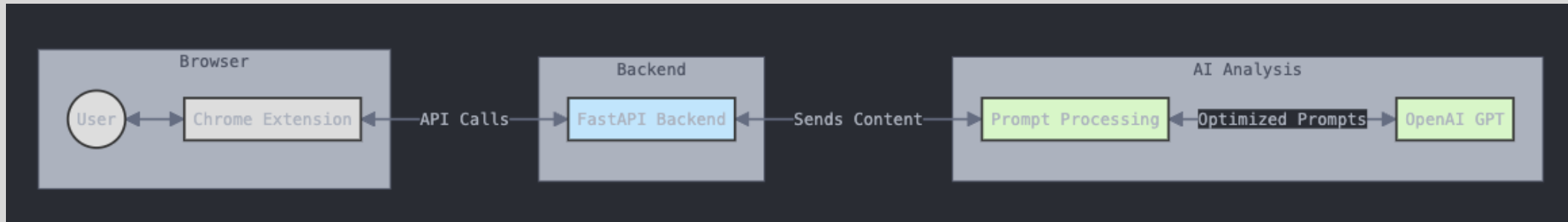
Matching Status



### Test Dataset

100 labeled samples for model evaluation

# Technical implementation



- Browser section: User interaction with Chrome Extension
- Backend: FastAPI server
- AI : prompt processing and LLM decision



# Prompt engineering

## Prompt Engineering



### Chain of Thought Example

#### Step-by-Step Analysis:

1. Understand the user's learning objective: '{label}'
2. Examine the page content: {content}
3. Identify key educational concepts in the content
4. Check if these concepts align with the user's learning goal
5. Consider whether this content would help achieve the learning objective
6. Evaluate if the content is worth the user's time and attention

#### Output:

Based on this analysis, output only 'True' if the content is relevant to the user's learning goal, or 'False' if it would be a distraction.

#### We included also:

- Few-shot learning
- Role prompting
- Zero-shot/Few-shot chain of thought
- Structured reasoning



# RESULTS

## Binary Classification Approach



True: Content Matches  
Goal

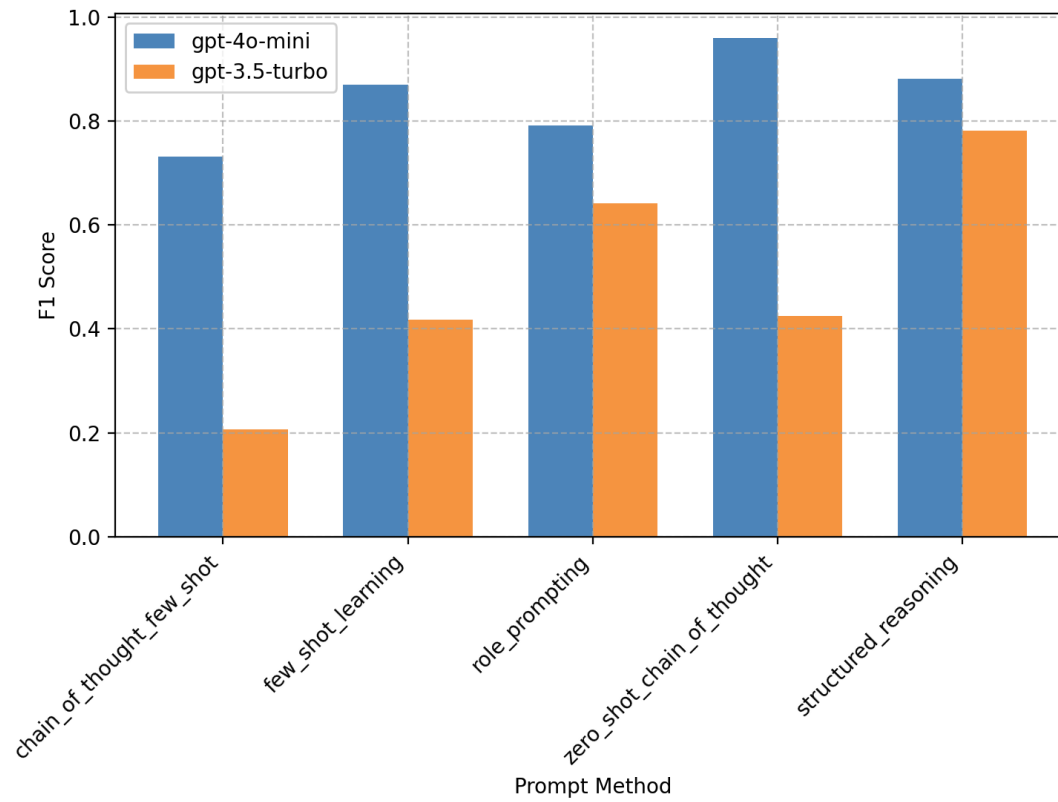


False: Content is  
Distraction

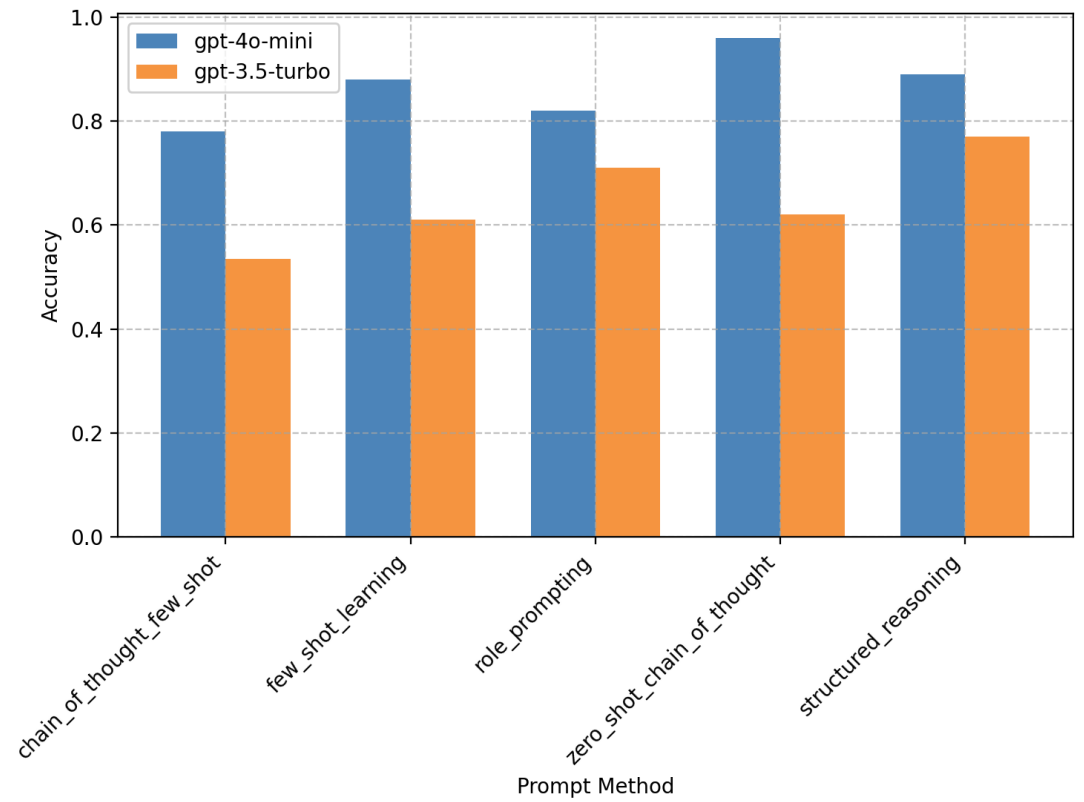
## Output Processing

Response validation: Extract "True" or "False" substring from LLM output

F1 Scores by Model and Prompt Method



Accuracy Scores by Model and Prompt Method





# PROMPT TUNING

## What is DSPy?

DSPy is a framework for programming foundation models that allows:

- Programming with LLMs using high-level primitives
- Automatic prompt optimization
- Systematic prompt engineering through modules
- Integration with various LLM providers (OpenAI, Anthropic, etc.)

## DSPy Optimizers

DSPy optimizers are tools that automatically improve prompts by:

- Learning from examples (bootstrapping)
- Testing different prompt variations
- Measuring performance using defined metrics
- Selecting the best performing prompts

Key optimizer: BootstrapFewShot - Automatically generates and tests few-shot examples

```
teleprompter = BootstrapFewShot(  
    metric=self.metric,  
    max_bootstrapped_demos=8,  
    max_labeled_demos=8,  
    max_rounds=10,  
)
```

# USING THE SAME DATASET

```
@dataclass
class Input:
    label: str
    content: str

@dataclass
class Output:
    classification: str
    reasoning: str = ""

class BrainrotDataset(Dataset):
    def __init__(self, data: pd.DataFrame):
        self.data = data
        # ...

class ZeroShotCoTClassifier(dspy.Module):
    # ...
```

## Code Structure - Classes

Key components in the code:

- Input/Output dataclasses - Define data structure
- BrainrotDataset - Custom dataset implementation
- ZeroShotCoTClassifier - Main classifier using chain-of-thought
- BrainrotOptimizer - Manages optimization process

# CLASS STRUCTURE

## ZeroShotCoTClassifier Implementation

The classifier uses step-by-step reasoning:

1. Analyzes user's learning goal
2. Evaluates content relevance
3. Assesses potential distractions
4. Makes binary classification (True/False)

Uses `dspy.Predict` for structured output generation

```
class ZeroShotCoTClassifier(dspy.Module):  
    def __init__(self):  
        super().__init__()  
        self.predictor = dspy.Predict("""Evaluate learning opportunity  
        Return True if content supports learning, False if it's a distr  
  
    def forward(self, label: str, content: str) -> Output:  
        steps = f"""Let's evaluate:  
        1. What is the user trying to learn? ({label})  
        2. What knowledge does this content provide?  
        3. Would it advance the user's goal?  
        4. Could it distract from the objective?  
        5. Is this the right time to engage?"""
```

# PROMPT USED



```
# Data splitting
train_val_size = int(0.8 * len(processed_data))
train_val_data = processed_data.iloc[:train_val_size]
test_data = processed_data.iloc[train_val_size:]

# Optimization
optimizer = BrainrotOptimizer(train_data, val_data)
results = optimizer.optimize_classifiers()

# Testing
best_name, best_classifier = optimizer.get_best_classifier()
# Final evaluation on test set
for _, row in test_data.iterrows():
    test_input = {
        'label': row['label'],
        'content': row['content']
    }
    prediction = best_classifier(**test_input)
```

## ✧ Optimization Process

The optimization workflow:

1. Data split into train/validation/test sets
2. BootstrapFewShot configured with parameters
3. Classifier compilation and optimization
4. Performance evaluation on validation set
5. Final testing on held-out test set

# SIMILARLY TO TRADITIONAL MODEL TRAINING



# INSIGHTS



## DSPy Optimizer Internals

Under the hood, DSPy optimizers work through:

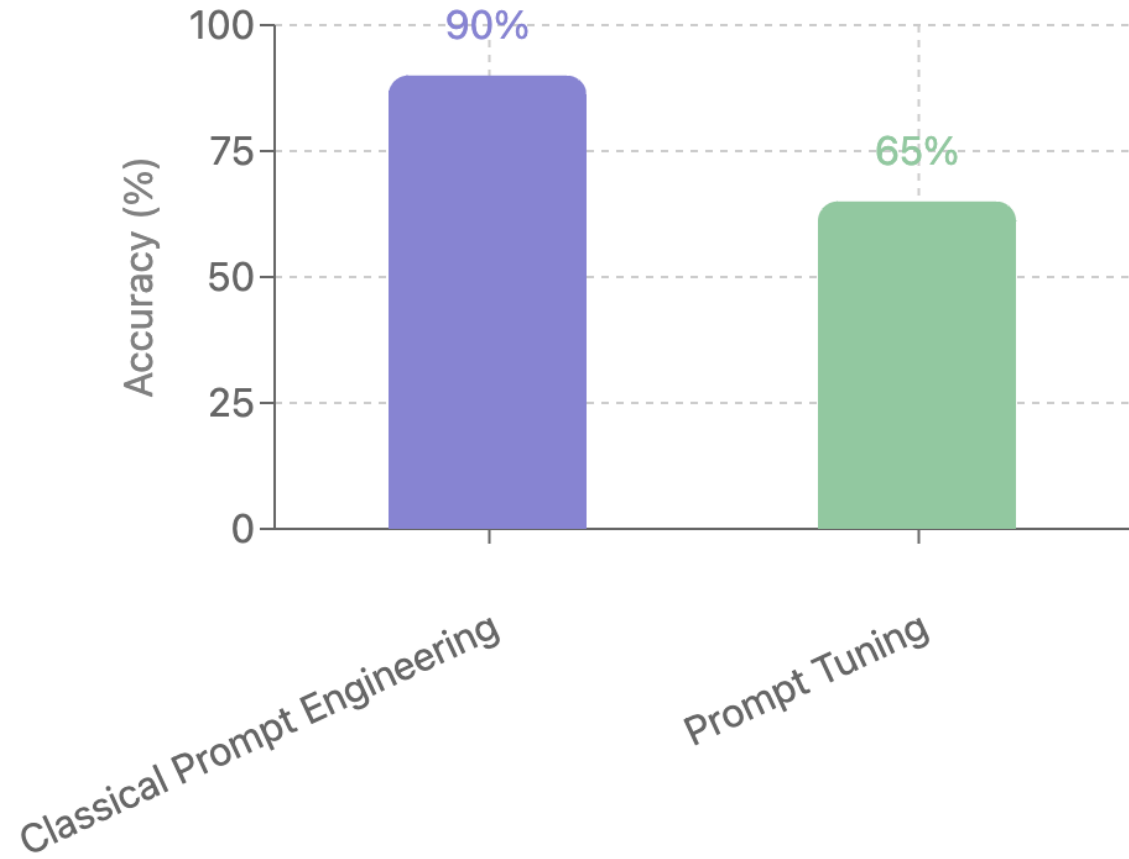
- LLM-driven Bootstrapping: Uses a 'critic' LLM (typically GPT-4) to evaluate and improve prompts
- Iterative Refinement: Conducts multiple rounds of testing and improvement (usually 5-10 rounds)
- Example Generation: Automatically creates new few-shot examples by analyzing successful and failed cases
- Prompt Evolution: Modifies instruction templates based on performance feedback
- Performance Tracking: Measures improvements using customizable metrics like accuracy or F1 score

# RESULTS

Why is it worse?

- Exact match for the output, instead of regex 'True/False'
- Library limitation in terms of compatibility of versions and dependencies, old documentation
- Too small dataset

## Zero Shot Chain of Thought Accuracy Comparison (GPT-4-mini)






# LangChain Ecosystem

Powerful tools for building AI applications





```
self.model = ChatOpenAI(  
    model_name=settings.OPENAI_MODEL,  
    api_key=settings.OPENAI_API_KEY  
)
```



# LangChain

- Framework for developing LLM applications
- Chains and prompts management
- Multiple LLM providers support
- Built-in memory systems
- Powerful agents and tools



```
os.environ["LANGCHAIN_API_KEY"] = settings.LANGCHAIN_API_KEY
os.environ["LANGCHAIN_TRACING_V2"] = settings.LANGCHAIN_TRACING_V2
os.environ["LANGCHAIN_PROJECT"] = settings.LANGCHAIN_PROJECT
```




# LangSmith

- Debug and monitor LLM applications
- Track costs and performance
- Evaluate model outputs
- Fine-tune prompts
- Collaborative development

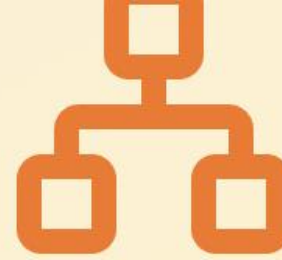
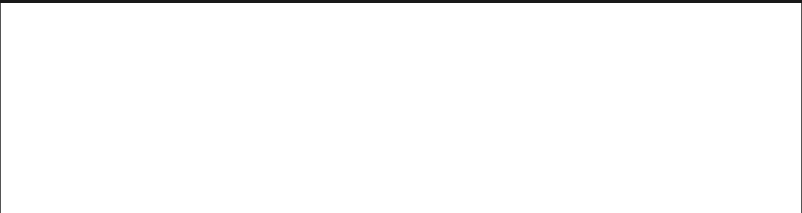








```
def get_workflow(self):  
    # Create workflow graph  
    workflow = StateGraph(MessagesState)  
  
    # Define nodes  
    workflow.add_node("agent", self.call_model)  
  
    # Set entry point  
    workflow.add_edge(START, "agent")  
  
    workflow.add_edge("agent", END)  
  
    return workflow
```



# LangGraph

- Build complex AI workflows
- Graph-based orchestration
- State management
- Parallel execution
- Advanced routing capabilities





# FOCUS AI GRAPH



**THANK YOU!**