# ILSP/AthenaRC submission: Constructing a Synthetic Multimodal Reasoning Dataset

## DCVLR 3rd place / IDC prize

Eleni E. Kasoura, Dimitris Damianos, **George Paraskevopoulos**, Vassilis Katsouros

**Institute for Language and Speech Processing, AthenaRC**

ATHENA'
Research & Innovation
Information Technologies

# Motivation

### Motivation

- Current work: We build LLMs for Greek (Meltemi, Llama-Krikri).

- Good quality instruction data are difficult to come by, especially for low-resource languages.

- Just Translating English datasets has limitations (culturally relevant LLMs, weird syntax, …)

### Aim

- Build a scalable, controllable and reproducible pipeline for image instruction synthesis with built-in quality checks.

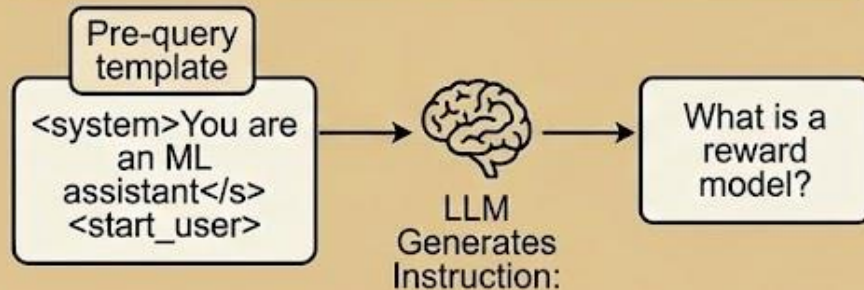- Systematically explore if we can heavily rely on synthetic data

# What is Magpie?

- Magpie (Xu et al., 2024): A self-synthesis framework for generating instruction–response pairs.

- **Scalable data synthesis**: Generates massive datasets directly from model weights, without human annotation.

- **Auto-Regressive Templates:** Inputs pre-query control tokens (e.g., <|start_header_id|>) to induce the spontaneous generation of diverse user queries and responses.

- **Good Downstream Performance:** Surpasses existing public datasets and rivals official proprietary models on benchmarks like AlpacaEval and ArenaHard.
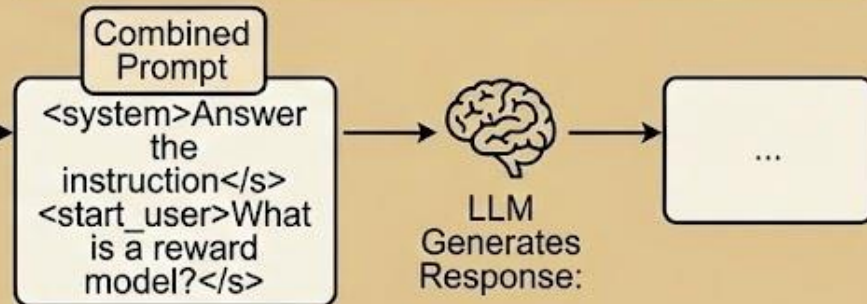
# Magpie Prompting

## Step 1: Instruction Generation

Pre-query template

`<system>You are an ML assistant</s> <start_user>`

→ LLM Generates Instruction: →

What is a reward model?

## Step 2: Response Generation

Combined Prompt

`<system>Answer the instruction</s> <start_user>What is a reward model?</s>`

→ LLM Generates Response: →

...

# Multimodal Magpie prompting



Step 1: Instruction Generation

Step 2: Response Generation

`<system>You are an ML assistant</s> <start_user> <visual_tokens>`

LLM Generates Instruction:

What is a reward model?

Combined Prompt

`<system>Answer the instruction</s> <start_user>What is a reward model?</s>`

LLM Generates Response:

...

# Pipeline

- Data selection and sampling (HuggingfaceM4/the_cauldron)
- Instruction pair generation
- Sanitization: Removes artifacts, ensures grounding to visual content.
- Optional rationale generation (reasoning) and tagging for complexity and domain.
- Reward Scoring and Filtering.

This process yields high-quality, diverse, vision-grounded instruction–response pairs.
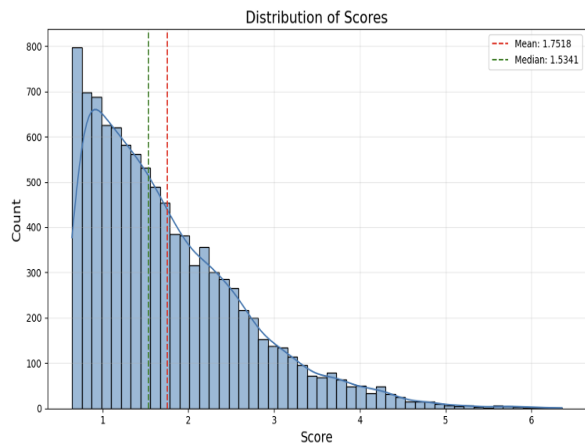
# System Prompt Categories



Figure 2: Reward score distribution of selected samples

**Prompt Types and Examples:**

- General – 'Ask a question about the objects in this image.'

- Math – 'Summarize the main statistical findings in this chart.'

- Spatial – 'Describe where the main objects are located.'

- Text – 'Read and summarize visible text.'

- Multi-image – 'Compare similarities or differences between two scenes.'

# Examples of different reward scoring QA

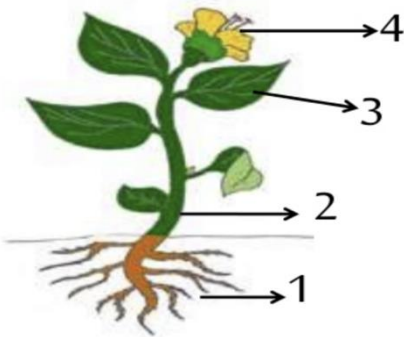*High reward scoring*

**"image_id"** : "ai2d_13",

**"system_prompt"** : "Create an instruction that asks the user to explain the relationship between variables or symbols in the image.",

**"instruction"** : "Explain how the roots, stems, leaves, and flowers of a plant are interconnected and their roles in the plant's life cycle.",

**"response"** : "The roots, stems, leaves, and flowers of a plant are interconnected and each plays a crucial role in the plant's life cycle. 1. **Roots**: The roots anchor the plant in the soil and absorb water and nutrients from the soil. They also store energy for the plant and help regulate the plant's water balance by absorbing water when it is available and releasing it when needed. Roots contain chlorophyll, which is essential for photosynthesis. 2. **Stems**: The stems support the plant and hold its leaves up so they can get sunlight. They transport water and nutrients absorbed by the roots throughout the plant. Stems also store energy and produce new leaves and flowers. 3. **Leaves**: Leaves are the primary site of photosynthesis, where the plant converts light energy into chemical energy (glucose) using carbon dioxide and water. Leaves also produce oxygen as a byproduct of photosynthesis. Leaves are attached to the stems and are responsible for the plant's growth and development. 4. **Flowers**: Flowers are the reproductive structures of the plant. They produce seeds, which contain the plant's genetic material. The flowers attract pollinators like bees and butterflies, which transfer pollen from one flower to another, facilitating fertilization. Once fertilized, the flowers develop into fruit, which contains the seeds. All these parts work together to ensure the plant's survival and reproduction. The roots provide the necessary nutrients and water, the stems support the plant and transport resources, the leaves perform photosynthesis, and the flowers facilitate reproduction. Together, they form a complex system that allows plants to grow, thrive, and reproduce.",

**"score"**: 4.701957702636719



Plants have 4 parts

# Examples of different reward scoring QA
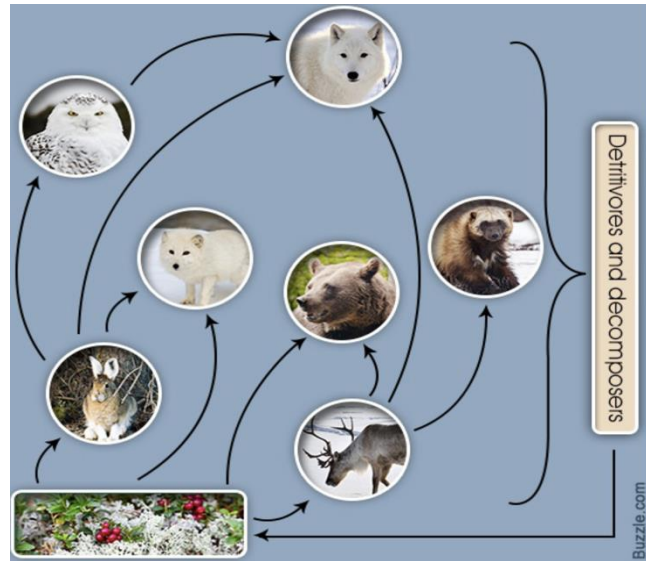
*Low reward scoring*

**"image_id"** : "ai2d_50",

**"system_prompt"** : "Ask a question about the arrangement of people, animals, or items in the image.",

**"instruction"** : "What is the difference between pinnate and palmate leaf arrangements as shown in the diagram?",

**"response"** : "The diagram illustrates two types of leaf arrangements: **pinnate** and **palmate**. 1. **Pinnate Leaf Arrangement (Sweet Clover, Bird's-foot Trefoil, Crownvetch):** - In a pinnate arrangement, the leaves are arranged along a central stem or rachis like the leaves of a feather. The can see this in the Sweet Clover, Bird's foot trefoil, and crownvetch. The leaves structure is that in pinnate the leaves are arranged along a a central stem, while in palmate the leaves are arranged at the end of the stem, resembling the hand of a fan. You can see this in the Kura clover, which has a palmate arrangement where all the leaflets radiate from a single point at the end of the petiole. In summary: - **Pinnate**: Leaves are arranged along a central stem. - **Palmate**: Leaves are arranged at the end of the stem.",

**"score"** : 0.7839621305465698

# Implementation Details

## Part 1

- Model: Qwen2.5-VL-7B-Instruct
- Dataset: Cauldron collection (≈50 datasets covering VQA, OCR, math, charts, etc.)
- Reward Model: OpenAssistant/reward-model-deberta-v3-large-v2
- Threshold: $t = 0$ (keep positive reward scores)
- Hardware: Leonardo-Booster Supercomputer, 4× NVIDIA A100 GPUs per node.

# Implementation Details

## Part 2

- System Prompts steer synthesis into task-specific directions:
  - Text extraction
  - Spatial reasoning
  - Mathematical reasoning
  - General comprehension

This controlled prompting is necessary, as it increases diversity and reduces bias toward simple image captions.

# Dataset Statistics

- Total: 10,000 curated samples.
- Average reward score: 1.75, median: 1.53.
- Dominant subset: Visual Question Answering (VQA).
- Underrepresented: Math diagrams and infographics.
- Selection favors clarity and instruction quality.

# Conclusion

- The proposed Magpie-style pipeline effectively synthesizes grounded visual instructions.

- Automated scoring ensures reproducibility and high data quality.

- System prompting is essential due to bias towards image captioning for diverse task coverage

# Limitations

- Limited ablations due to time/compute constraints.
- Qwen2.5-VL-7B is a relatively small model (7B parameters).
- Reward filtering strategy can be improved.
- Expand task diversity with new system prompts.

# DCVLR Competition



Final Leaderboard

DCVLR 1st Edition - NeurIPS 2025 Competition Results

| Rank | Team | Dataset | Technical Report | Accuracy | Δ vs Baseline |
|---|---|---|---|---|---|
| 1 🏅 | University of Hawaii at Manoa - team 3 | View Dataset | View Report | 46.12% | +7.28% |
| 2 🥈 | AFIE | View Dataset | View Report | 42.13% | +3.29% |
| 3 🥉 | AthenaRC | View Dataset | View Report | 41.81% | +2.98% |

https://dcvlr-neurips.github.io/

# AIXPERT

## What is AIXPERT

AIXPERT is an international research initiative funded by the European Union's Horizon Europe programme and the Swiss State Secretariat for Education, Research and Innovation (SERI). Our mission is to make AI smarter, safer, and more trustworthy across critical sectors such as healthcare, human resources, manufacturing, robotics, and the creative industries.

- Build an adaptable, explainable AI-agentic platform

- Define and assess AI trustworthiness
- Advance explainable multimodal foundation models
- Demonstrate real-world impact through pilot use cases (healthcare, recruitment, educational robotics, manufacturing and creative arts)

Έρευνα & Καινοτομία
Τεχνολογίες Πληροφορίας **ΑΘΗΝΑ'** **ATHENA'** Research & Innovation
Information Technologies

# Thank you!