Towards Resource-Efficient Compound AI Systems

Gohar Irfan Chaudhry^{*} Esha Choukse[†] Íñigo Goiri[†] Rodrigo Fonseca[†] Ricardo Bianchini[‡] Adam Belay^{*}

MIT CSAIL * Microsoft Azure Research – Systems † Microsoft Azure ‡













Are these systems efficient?















Video Q/A: workflow-specific knobs



Video Q/A: num. frames



Video Q/A: STT



Video Q/A: LLM



Video Q/A: hardware configuration



Video Q/A: parallelism





Outline

1. Compound AI workflow: Video Q/A

- a. Possible configurations.
- b. Workflow implementation.
- c. Workflow execution.
- 2. Our vision.
- 3. Murakkab.

Video Q/A: configuring with only 2 knobs



Video Q/A: manual configuration



Video Q/A: optimized configuration











Manual workflow configuration is suboptimal.

Outline

1. Compound AI workflow: Video Q/A

- a. Possible configurations.
- b. Workflow implementation.
- c. Workflow execution.
- 2. Our vision.
- 3. Murakkab.


```
/* ===== Nodes of the workflow (simplified) ===== */
scene_detection = Tool(name="OpenCV", key=AWS_SSH_KEY,
                       params={num frames: 15},
                       resources={CPUs: 32})
speech to text = MLModel(name="Whisper", key=OPENAI API KEY,
                          resources={PTUs: 50})
object_detection = MLModel(name="CLIP", key=AZURE_SSH_KEY,
                          resources={CPUs: 2})
question_answer = LLM(name="Llama", key=DATABRICKS_API_KEY,
                      params={batch size: 256},
                      resources={GPUs: 4,GPU_Type: H100},
                      system_prompt="You are an agent that can understand videos.",
                      user_prompt="Answer the given question about the video.")
/* ===== 0uerv and data ===== */
question = "What is the name of the person wearing the red dress?"
videos = ["road trip.mp4"]
/* ===== Workflow (ordering of nodes/data flow) ===== */
audio, frames = scene detection(videos)
transcript = speech_to_text(audio)
refined_frames = object_detection(frames)
              = question_answer(question, transcript, frames)
answer
```


Models/Tools

```
/* ===== Nodes of the workflow (simplified) ===== */
scene_detection = Tool(name="OpenCV", key=AWS_SSH_KEY,
                       params={num frames: 15},
                       resources={CPUs: 32})
speech_to_text = MLModel(name="Whisper", key=0PENAI_API_KEY,
                          resources={PTUs: 50})
object_detection = MLModel(name="CLIP", key=AZURE_SSH_KEY,
                          resources={CPUs: 2})
question_answer = LLM(name="Llama", key=DATABRICKS_API_KEY,
                      params={batch size: 256},
                      resources={GPUs: 4,GPU_Type: H100},
                      system_prompt="You are an agent that can understand videos.",
                      user_prompt="Answer the given question about the video.")
/* ===== Query and data ===== */
question = "What is the name of the person wearing the red dress?"
videos = ["road trip.mp4"]
/* ===== Workflow (ordering of nodes/data flow) ===== */
audio, frames = scene detection(videos)
transcript = speech_to_text(audio)
refined_frames = object_detection(frames)
              = question_answer(question, transcript, frames)
answer
```


Models/Tools

```
/* ===== Nodes of the workflow (simplified) ===== */
scene_detection = Tool(name="OpenCV", key=AWS_SSH_KEY,
                       params={num frames: 15},
                       resources={CPUs: 32})
speech_to_text = MLModel(name="Whisper", key=OPENAI_API_KEY,
                          resources={PTUs: 50})
object_detection = MLModel(name="CLIP", key=AZURE_SSH_KEY,
                          resources={CPUs: 2})
question_answer = LLM(name="Llama", key=DATABRICKS_API_KEY,
                      params={batch size: 256},
                      resources={GPUs: 4,GPU_Type: H100},
                      system prompt="You are an agent that can understand videos.",
                      user_prompt="Answer the given question about the video.")
/* ===== Query and data ===== */
question = "What is the name of the person wearing the red dress?"
videos = ["road trip.mp4"]
/* ===== Workflow (ordering of nodes/data flow) ===== */
audio, frames = scene detection(videos)
transcript = speech_to_text(audio)
refined_frames = object_detection(frames)
              = question_answer(question, transcript, frames)
answer
```

Task-specific parameters

	<pre>/* ===== Nodes of the workflow (simplified) ===== */ scene_detection = Tool(name="OpenCV", key=AWS_SSH_KEY,</pre>
Models/Tools	<pre>speech_to_text = MLModel(name="Whisper", key=OPENAI_API_KEY,</pre>
	<pre>object_detection = MLModel(name="CLIP", key=AZURE_SSH_KEY,</pre>
Resources	<pre>question_answer = LLM(name="Llama", key=DATABRICKS_API_KEY,</pre>
	<pre>/* ===== Query and data ===== */ question = "What is the name of the person wearing the red dress?" videos = ["road_trip.mp4"]</pre>
	<pre>/* ===== Workflow (ordering of nodes/data flow) ===== */ audio, frames = scene_detection(videos) transcript = speech_to_text(audio) refined_frames = object_detection(frames) answer = question_answer(question, transcript, frames)</pre>

Task-specific parameters

	<pre>/* ===== Nodes of the workflow (simplified) ===== */ scene_detection = Tool(name="OpenCV", key=AWS_SSH_KEY,</pre>	Task-specific
Models/Tools	<pre>speech_to_text = MLModel(name="Whisper", key=OPENAI_API_KEY,</pre>	parameters
	<pre>object_detection = MLModel(name="CLIP", key=AZURE_SSH_KEY,</pre>	Multiple
Resources	<pre>question_answer = LLM(name="Llama", key=DATABRICKS_API_KEY,</pre>	providers
	<pre>/* ===== Query and data ===== */ question = "What is the name of the person wearing the red dress?" videos = ["road_trip.mp4"]</pre>	
	<pre>/* ===== Workflow (ordering of nodes/data flow) ===== */ audio, frames = scene_detection(videos) transcript = speech_to_text(audio) refined_frames = object_detection(frames) answer = question_answer(question, transcript, frames)</pre>	

	<pre>/* ===== Nodes of the workflow (simplified) ===== */ scene_detection = Tool(name="OpenCV", key=AWS_SSH_KEY,</pre>	Task-specific
Models/Tools	<pre>speech_to_text = MLModel(name="Whisper", key=OPENAI_API_KEY,</pre>	parameters
	<pre>object_detection = MLModel(name="CLIP", key=AZURE_SSH_KEY,</pre>	Multiple
Resources	<pre>question_answer = LLM(name="Llama", key=DATABRICKS_API_KEY,</pre>	providers
	/* ===== Query and data ===== */ question = "What is the name of the person wearing the red dress?" videos = ["road_trip.mp4"]	Application logic
	<pre>/* ===== Workflow (ordering of nodes/data flow) ===== */ audio, frames = scene_detection(videos) transcript = speech_to_text(audio) refined_frames = object_detection(frames) answer = question_answer(question, transcript, frames)</pre>	

	<pre>/* ===== Nodes of the workflow (simplified) ===== */ scene_detection = Tool(name="OpenCV", key=AWS_SSH_KEY,</pre>	Task-specific parameters
Models/Tools	<pre>speech_to_text = MLModel(name="Whisper", key=OPENAI_API_KEY,</pre>	
	<pre>object_detection = MLModel(name="CLIP", key=AZURE_SSH_KEY,</pre>	Multiple
Resources	<pre>question_answer = LLM(name="Llama", key=DATABRICKS_API_KEY,</pre>	providers
	<pre>/* ===== Query and data ===== */ question = "What is the name of the person wearing the red dress?" videos = ["road_trip.mp4"]</pre>	Application logic

Tightly coupled configuration and application logic.
Outline

1. Compound AI workflow: Video Q/A

- a. Possible configurations.
- b. Workflow implementation.
- c. Workflow execution.
- 2. Our vision.
- 3. Murakkab.

```
Workflow(
    models=[A,B,C],
    api_keys=[foo,bar],
    stages=[S1,S2,S3],
    knobs={x:5,y:True}
)
```

Explicitly defined workflow logic coupled with configuration details.



Explicitly defined workflow logic coupled with configuration details.

Static DAG with fixed model(s) and other knobs.



coupled with configuration details.

Static DAG with fixed model(s) and other knobs.



Models and resources managed by different entities/providers.



Models and resources managed by different entities/providers.





Models and resources managed by different entities/providers.



Fragmented workflows spanning different entities with limited inter-layer visibility.



Are these systems efficient?

Are these systems efficient?





Achieving high efficiency while satisfying SLOs is difficult in today's compound AI systems.

Why?

Why?

1. Rigid workflows:

a. Tight coupling between configurations and application logic.

•••
<pre>/* ===== Nodes of the workflow (simplified) ===== */ scene_detection = Tool(name="OpenCV", key=AWS_SSH_KEY,</pre>
<pre>speech_to_text = MLModel(name="Whisper", key=OPENAI_API_KEY,</pre>
<pre>object_detection = MLModel(name="CLIP", key=AZURE_SSH_KEY,</pre>
<pre>question_answer = LLM(name="Llama", key=DATABRICKS_API_KEY,</pre>
<pre>/* ===== Query and data ===== */ question = "What is the name of the person wearing the red dress?" videos = ["road_trip.mp4"]</pre>
<pre>/* ===== Workflow (ordering of nodes/data flow) ===== */ audio, frames = scene_detection(videos) transcript = speech_to_text(audio) refined_frames = object_detection(frames) accumention_accument(superior_transcript_frames)</pre>
answer – quesceon_answer(quesceon; craitser(pc; rraites)

Why?

1. Rigid workflows:

a. Tight coupling between configurations and application logic.

2. Limited cross-layer visibility:

- a. Workflow orchestration and resource management.
- b. Multiple model/tool providers.





Outline

1. Compound Al workflow: Video Q/A

- a. Possible configurations.
- b. Workflow implementation.
- c. Workflow execution.

2. Our vision.

3. Murakkab.



Poor efficiency!

High efficiency!

• High level <u>descriptions</u> (in natural language):

ob_description	= "What is the name of	the person	wearing the	red dress?"
nputs	<pre>= ["road_trip.mp4"]</pre>	·	-	
result	<pre>= run(job_description,</pre>	inputs,)	

• High level <u>descriptions</u> (in natural language) and <u>SLOs</u>:

job_description inputs 5LOs	<pre>= "What is the name of = ["road_trip.mp4"] = {cost: "<\$2"}</pre>	the person wearing the red dress?"
result	<pre>= run(job_description,</pre>	inputs, SLOs)

• High level <u>descriptions</u> (in natural language) and <u>SLOs</u>:



• High level <u>descriptions</u> (in natural language) and <u>SLOs</u>:



Decouple configuration from application logic.

Idea 2: unified compound AI system

Idea 2: unified compound AI system

- <u>Single entity</u> manages:
 - Workflow orchestration
 - Models/tools and resources

Idea 2: unified compound AI system

• <u>Single entity</u> manages:

- Workflow orchestration
- Models/tools and resources





Outline

1. Compound Al workflow: Video Q/A

- a. Possible configurations.
- b. Workflow implementation.
- c. Workflow execution.
- 2. Our vision.
- 3. Murakkab.

Declarative Job Description + SLOs











Agents and Tools Library

(LLMs, ML models, tools) [e.g., DeepSeek-R1, Whisper, Web Browser API]



[e.g., DeepSeek-R1, Whisper, Web Browser API]

Runtime System + Infrastructure

(Resource monitoring, allocation, scheduling, auto-scaling)



Murakka	ab					
 <u>SLOs:</u> 1. Price (\$) 2. Quality (% or tier) 3. Latency (s) 	Efficiency: 1. Cost (\$) 2. Energy (kWh) 3. Resources (e.g., # GPUs)					
What is the name of the person wearing the red dress?" [SLO = < \$2]	LLM planner	Placeholo	Her Nodes	Model-selection +	Model Instances	
				WORKTIOW-KNODS	(e.g., Whisper on GPUs 1,2 Gemma-3 on GPU 3)	
Agents and Tools Library (LLMs, ML models, tools) [e.g., DeepSeek-R1, Whisper, Web Browser API]		Runtim (Re	ne System + source monitorin scheduling, auto	Infrastructure g, allocation, o-scaling) 67		

Efficiency:

Efficiency:

• Reallocate models and resources with changing workflows/load.

Efficiency:

- Reallocate models and resources with changing workflows/load.
- Collocation and model sharing:
 - More batching reduces marginal cost of serving new workflows.

Efficiency:

- Reallocate models and resources with changing workflows/load.
- Collocation and model sharing:
 - More batching reduces marginal cost of serving new workflows.

Quality:

• Open-source datasets/benchmarks as a proxy for new queries.

Efficiency:

- Reallocate models and resources with changing workflows/load.
- Collocation and model sharing:
 - More batching reduces marginal cost of serving new workflows.

- Open-source datasets/benchmarks as a proxy for new queries.
- Periodic feedback from users.
Video Q/A: optimized configuration



Dynamic workflow optimization





Are these systems efficient? NO! But they can be...

Are these systems efficient? NO! But they can be...

Decoupling app. logic from configuration
Unified compound AI system

Efficiency:

Quality:

Efficiency:

- Different hardware configurations.
- Multi-workflow joint optimization.

Quality:

Efficiency:

- Different hardware configurations.
- Multi-workflow joint optimization.

Quality:

- Per-task quality vs end-to-end quality.
- Closed source models (e.g., GPT-4o).

Efficiency:

- Different hardware configurations.
- Multi-workflow joint optimization.

Quality:

- Per-task quality vs end-to-end quality.
- Closed source models (e.g., GPT-4o).

End-to-end automatic profiling and optimization framework for multiple workflows on shared infrastructure.

Outline

- 1. Compound AI workflow: Video Q/A
 - a. Possible configurations.
 - b. Workflow implementation.
 - c. Workflow execution.
- 2. Our vision.
- 3. Murakkab.

girfan@mit.edu