# Tracking the Evolution of R and Python Tools for GenAI:
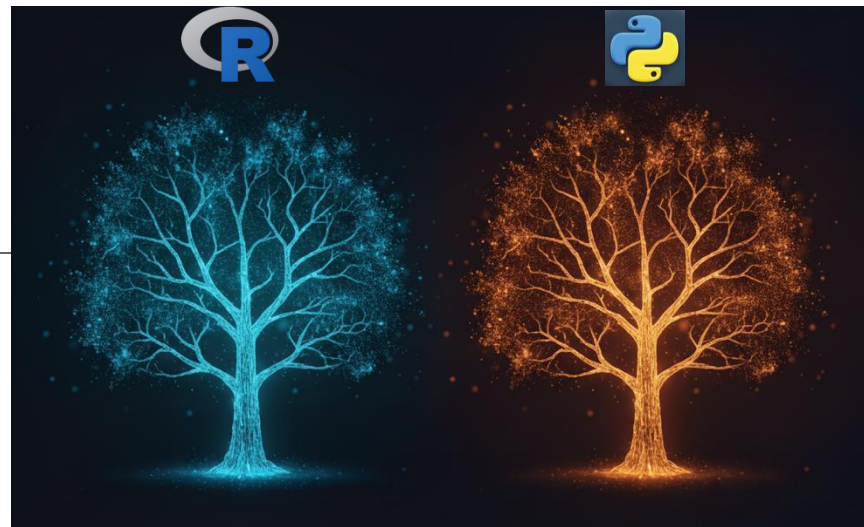
## State of the Ecosystem & Lessons from Building in Both

Xu Fei
Senior Solutions Engineer
xu@a2-ai.com
linkedin.com/in/xufei1/

xuf12

# Who is this talk for

- Language: R-primary or multilingual teams

- Background: Have experience building apps, interested in adding GenAI features (e.g. chatbot, coding assistant)

A2-Ai

# Our experience – chatbot for internal tools

ChatGPT
(generic chatbot)

How do I install dplyr?

To install the `dplyr` package in R, you can use one of these methods depending on where you're working (R console, RStudio, etc.):

💻 **Standard CRAN Installation**

Open your R console or RStudio and run:

```r
install.packages("dplyr")
```
⊡ Copy code

A2-Ai

# Our experience – chatbot for internal tools

### ChatGPT
### (generic chatbot)

How do I install dplyr?

To install the `dplyr` package in R, you can use one of these methods depending on where you're working (R console, RStudio, etc.):

💻 **Standard CRAN Installation**

Open your R console or RStudio and run:

```r
install.packages("dplyr")
```
📋 Copy code

### PharmaSage
### (uses internal docs + code)

How do I install dplyr?

**Installing dplyr with RV**

**Primary Solution: Using RStudio (Recommended)**

In RStudio:

1. Open your project in RStudio
2. First verify RV is active by running in the R console:

```
.rv$summary()
```

A2-Ai

# Our experience – chatbot with tool calling

# Our experience – MCP server built in R
## Enables natural language queries to do exploratory data analysis



> Great thanks, can you update the legend label to instead say Treatment group and dofetilide and placebo instead of the doses?

# 3 core building blocks of GenAI applications

1. Model Interface

      - Managing LLM API calls

Model Interface

# 3 core building blocks of GenAI applications

1. Model Interface

    - Managing LLM API calls

2. Context Handling

    - Handling anything that goes into the LLM

Context Handling

Model Interface

# 3 core building blocks of GenAI applications

1. Model Interface
   - Managing LLM API calls

2. Context Handling
   - Handling anything that goes into the LLM

3. Evaluation
   - Validating the results

| Context Handling | Evaluation |
| :---: | :---: |

| Model Interface |
| :---: |

# 3 core building blocks of GenAI applications

1. Model Interface
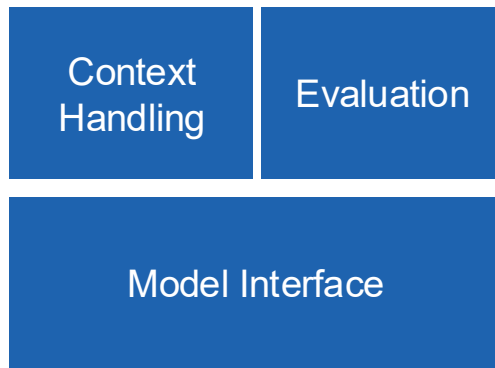  - Managing LLM API calls

2. Context Handling
  - Handling anything that goes into the LLM

3. Evaluation
  - Validating the results

| Context Handling | Evaluation |
|---|---|
| Model Interface | |

# Model Interface
You might like to ask an LLM a question in this way...

```
chat <- chat_anthropic(model = "claude-haiku-4-5-20251001")
chat$chat("How do I install dplyr?")
```

# Model Interface

...but sending this question in raw HTTP actually looks like this

```r
# a bare minimum LLM call
library(httr2)
response <- request("https://api.anthropic.com/v1/messages") |>
  req_headers(
    "x-api-key" = Sys.getenv("ANTHROPIC_API_KEY"),
    "anthropic-version" = "2023-06-01"
  ) |>
  req_body_json(list(
    model = "claude-haiku-4-5-20251001",
    max_tokens = 1024,
    messages = list(list(role = "user", content = "How do I install dplyr?"))
  )) |>
  req_perform() |>
  resp_body_json()
```

A2-Ai

# Model Interface – client SDKs

# Model Interface – SDK

Anthropic

**USING THE API**

## Client SDKs

Python

TypeScript

Java

Go

C#

Ruby

PHP

Beta namespace in client SDKs

← Where is R?

# Model Interface – SDK

## OpenAI

```
1   from openai import OpenA
2
3   client = OpenAI()
4
5   response = client.respons
6     model="gpt-4.1",
7     input="Tell me a three sentence bedtime st
8   )
9
10  print(response)
```

curl

javascript

✓ **python**

csharp

Also no R…

A2-Ai

# Model Interface – SDK

Amazon Bedrock

Documentation  >  Amazon Bedrock  >  User Guide

## Hello Amazon Bedrock

⤓ PDF    ⤓ RSS    ⬤ Focus mode

The following code examples show how to get started using Amazon Bedrock.

| .NET | Go | JavaScript | **Python** | Swift |
|------|-----|-----------|------------|-------|

Still no R ☹

# Model Interface – SDK

Amazon Bedrock



## Hello Amazon Bedrock

[download] PDF    [download] RSS    ⊙ Focus mode

The following code examples show how to get started using Amazon Bedrock.

| .NET | Go | JavaScript | **Python** | Swift |
|------|-----|-----------|-----------|-------|

but


Paws

A2-Ai

# Model Interface – pick a direction?

# Model Interface – the answer for R



Supports Anthropic, OpenAI, AWS Bedrock, Gemini, etc.

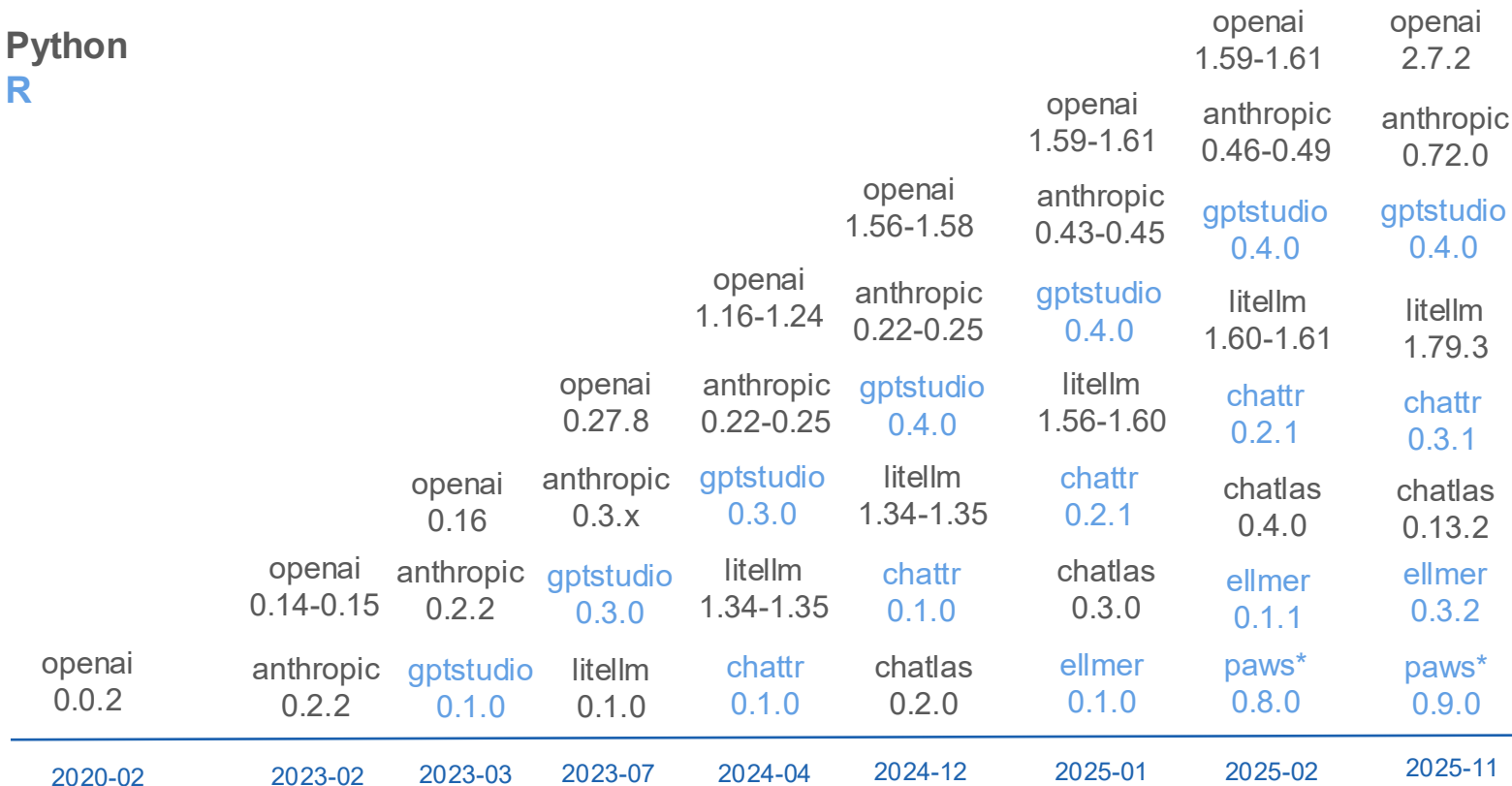0.3.2 (latest) released on 2025-09-03
0.1.0 released on 2025-01-09

# Model Interface – ellmer

```r
library(ellmer)
chat <- chat_anthropic(model = "claude-haiku-4-5-20251001")
chat$chat("How do I install dplyr?")
```

# Model Interface – ecosystem evolution

**Python**
**R**

| 2020-02 | 2023-02 | 2023-03 | 2023-07 | 2024-04 | 2024-12 | 2025-01 | 2025-02 | 2025-11 |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | openai 1.59-1.61 | openai 2.7.2 |
| | | | | | | openai 1.59-1.61 | anthropic 0.46-0.49 | anthropic 0.72.0 |
| | | | | | openai 1.56-1.58 | anthropic 0.43-0.45 | gptstudio 0.4.0 | gptstudio 0.4.0 |
| | | | | openai 1.16-1.24 | anthropic 0.22-0.25 | gptstudio 0.4.0 | litellm 1.60-1.61 | litellm 1.79.3 |
| | | | openai 0.27.8 | anthropic 0.22-0.25 | litellm 1.56-1.60 | chattr 0.2.1 | chattr 0.3.1 |
| | | openai 0.16 | anthropic 0.3.x | gptstudio 0.3.0 | litellm 1.34-1.35 | chattr 0.2.1 | chatlas 0.4.0 | chatlas 0.13.2 |
| | openai 0.14-0.15 | anthropic 0.2.2 | gptstudio 0.3.0 | litellm 1.34-1.35 | chattr 0.1.0 | chatlas 0.3.0 | ellmer 0.1.1 | ellmer 0.3.2 |
| openai 0.0.2 | anthropic 0.2.2 | gptstudio 0.1.0 | litellm 0.1.0 | chattr 0.1.0 | chatlas 0.2.0 | ellmer 0.1.0 | paws* 0.8.0 | paws* 0.9.0 |

\* Not only for LLM model interface

# Model Interface – abstraction level

|  | R | Python |
|---|---|---|
| High-level orchestration frameworks | ? | langchain<br>llamaindex<br>autogen<br>etc. |
| Unified LLM abstraction | ellmer<br>gptstudio<br>chattr | litellm<br>chatlas<br>etc. |
| LLM provider SDKs | paws* | openai<br>anthropic<br>gemini<br>etc. |

| LLM provider REST API |
|---|

* Not only for LLM model interface

# Model Interface – # of monthly downloads

**Python**
**R**

|  | monthly downloads across all versions |
|---|---|
| openai | 128M |
| litellm | 39M |
| anthropic | 29M |
| chatlas | 16.5k |
| ellmer | 9.2k |
| paws* | 4.1k |
| chattr | 1.5k |
| gptstudio | 1.1k |

Numbers calculated using R library "cranlogs" and pypistats.org as of November 10, 2025
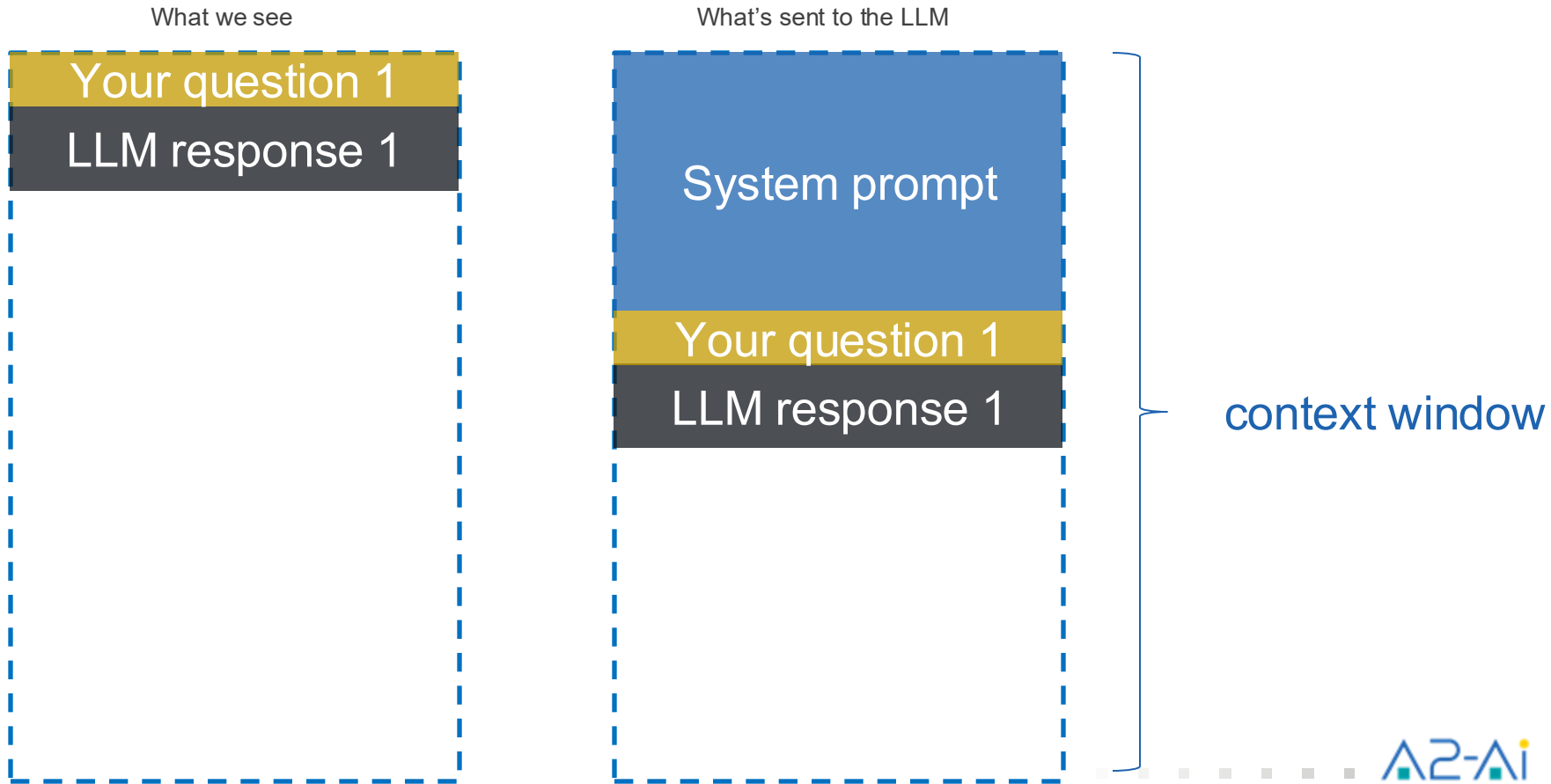* Not only for LLM model interface

# Model Interface – case study
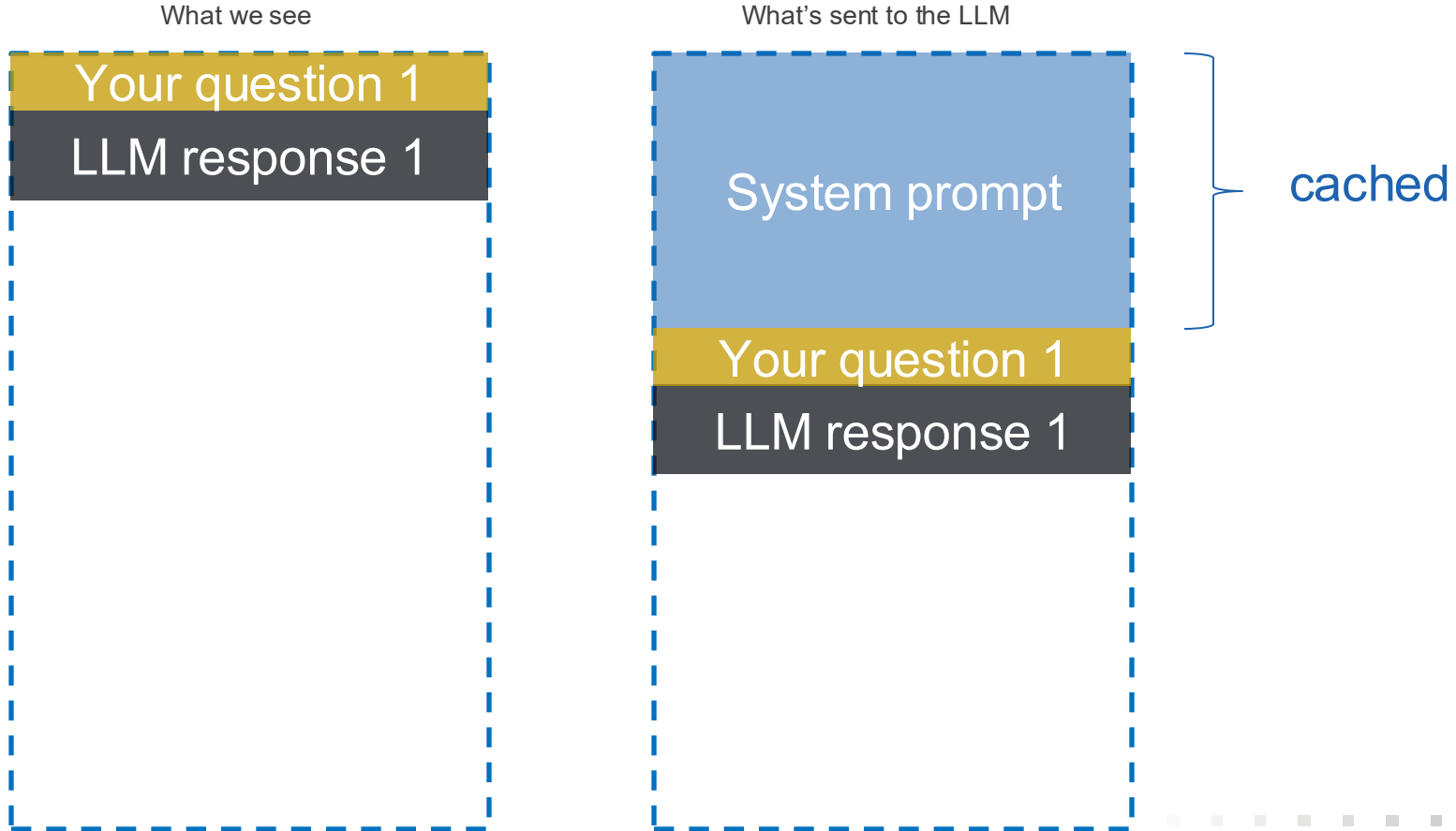## Able to extend model interface using Python SDK

Need:

We need prompt caching to reduce costs in AWS Bedrock.

A2-Ai

# Context handling: what's in a context window?

**What we see**

| |
|---|
| Your question 1 |
| LLM response 1 |

**What's sent to the LLM**

| |
|---|
| System prompt |
| Your question 1 |
| LLM response 1 |

context window

# Context handling: what's cached?

What we see

What's sent to the LLM

Your question 1

LLM response 1

System prompt

Your question 1

LLM response 1

cached

# Model Interface – case study
## Able to extend model interface using Python SDK

**Need:**

We need prompt caching to reduce costs in AWS Bedrock.

**Challenge:**

In September 2025, it was not part of the chatlas package for the Anthropic models that we used.

A2-Ai

# Model Interface – case study
## Able to extend model interface using Python SDK

Need:

   We need prompt caching to reduce costs in AWS Bedrock.

Challenge:

   In September 2025, it was not part of the chatlas package for the Anthropic
   models that we used.

Solution:

   We decided to add it to chatlas ourselves, because prompt caching was in
   the AWS Bedrock Python SDK.

A2-Ai

# Model Interface – REST API
## Without the SDK, we would have to deal with raw HTTP directly

```r
# a bare minimum LLM call
library(httr2)
response <- request("https://api.anthropic.com/v1/messages") |>
  req_headers(
    "x-api-key" = Sys.getenv("ANTHROPIC_API_KEY"),
    "anthropic-version" = "2023-06-01"
  ) |>
  req_body_json(list(
    model = "claude-haiku-4-5-20251001",
    max_tokens = 1024,
    messages = list(list(role = "user", content = "How do I install dplyr?"))
  )) |>
  req_perform() |>
  resp_body_json()
```

A2-Ai

# Model Interface – adding a new feature
## Able to extend model interface using Python SDK



```
chatlas/_provider_anthropic.py                                    +147 -6 ████  ···

398  +    def _get_system_content(self, system_text: str):
399  +        """Get system content with optional caching."""
400  +        if self._cache_system_prompt:
401  +            # Build cache control with optional TTL
402  +            cache_control = {"type": "ephemeral"}
403  +            if self._cache_ttl:
404  +                cache_control["ttl"] = self._cache_ttl
405  +
406  +            # Return system content as a list with
         cache_control
407  +            return [
408  +                {
409  +                    "type": "text",
410  +                    "text": system_text,
411  +                    "cache_control": cache_control
412  +                }
```

xuf12 committed on Sep 24

**Add prompt caching for Anthropic (including Bedrock) models**

⎇ add-prompt-caching-bedrock-anthropic

```
chatlas/_provider_anthropic.py                                    +147 -6 █████  ···

@@ -74,6 +74,8 @@ def ChatAnthropic(
74       api_key: Optional[str] = None,              74       api_key: Optional[str] = None,
75       max_tokens: int = 4096,                     75       max_tokens: int = 4096,
76       kwargs: Optional["ChatClientArgs"] = None,  76       kwargs: Optional["ChatClientArgs"] = None,
                                                      77  +    cache_system_prompt: bool = False,
                                                      78  +    cache_ttl: Optional[Literal["5m", "1h"]] = None,
77       ) -> Chat["SubmitInputArgs", Message]:      79       ) -> Chat["SubmitInputArgs", Message]:
```

# 3 core building blocks of GenAI applications
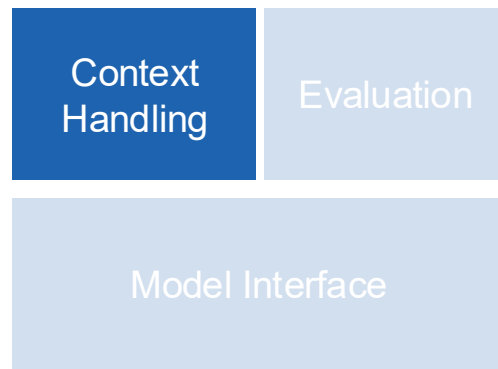
1. Model Interface

- Managing LLM API calls

## 2. Context Handling

- Handling anything that goes into the LLM

3. Evaluation

- Validating the results are as intended

| Context Handling | Evaluation |
| --- | --- |
| Model Interface | |

A2-Ai

# Context handling
## Example for a knowledge assistant app

← Submitting questions via model interface

all of the work to process
enterprise-specific documents and
code to share with the model

(requires additional tooling)

A2-Ai

# Context handling – pick a direction again?

# Context handling

1. In R, start with ellmer and btw

# Context handling
## btw can use your R session as context



EXPLORER ···

-AND-AI-TALK-PREP

- .gitignore       U
- .Renviron        U
- .Rprofile        U
- llm-calls.qmd    U
- README.md        U

llm-calls.qmd ⓤ ✕ · {} ## Context Untitle

🖥 Preview ☐ Render on Save · Source Visua

llm-calls.qmd > ...

```
30
   ▷ Run Cell | Run Next Cell | Run Above
31  ```{r}
32  library(cranlogs)
33  pkg <- "ellmer"
34  # last-month is last 30 days
    per API; you can also use
    explicit dates
35  df <- cran_downloads(packages
    = pkg, when = "last-month")
36  monthly_total <- sum(df$count)
37  df; cat("Monthly total:",
    monthly_total, "\n")
```

SESSION   CONNECTIONS   HELP   VIEWER   CHAT

← →   http://127.0.0.1:7957/

> Chat with {btw} tools

what can you see in total?

Looking at your R environment and project files, here's everything I can see:

## R Environment Objects (3 items):

1. `df` - A data frame with 30 rows of ellmer package download data:

   - Dates: 2025-10-11 to 2025-11-09
   - Daily download counts: ranging from 125 to 556
   - All entries for package "ellmer"

2. `monthly_total` - A numeric value of 9,259 (total downloads)

3. `pkg` - A character string containing "ellmer"

## Project Files:

- **README.md** (20 bytes) - Very small README

# Context handling
## Options when required libraries don't exist in R

1. In R, start with ellmer and btw

2. When required libraries don't exist in R, try to use reticulate to call specific Python libraries and functions

# Context handling
## Options when required libraries don't exist in R

**Key Steps**

### 1. Document Processing

`ragnar` works with a wide variety of document types, using [MarkItDown](#) to convert content to Markdown.

**ragnar / R / ragnar-package.R**

```
 9  ∨    .onLoad <- function(libname, pkgname) {
10        Sys.setenv(RETICULATE_PYTHON = "managed")
11        S7::methods_register()
12        reticulate::py_require(c(
13          # Pin onnxruntime until this is resolved:
14          #  https://github.com/microsoft/markitdown/issues/1266
15          # New VC++ version requirement begins:
16          # https://github.com/Microsoft/onnxruntime/releases/tag/v1.21.0
17          if (is_windows()) "onnxruntime<=1.20.1",
18          "markitdown[all]"
19        ))
```

# Context handling
## For multilingual teams

1. In R, start with ellmer and btw

2. When required libraries don't exist in R, try to use reticulate to call specific Python libraries and functions

3*. Create REST APIs for existing R (plumber) or Python (FastAPI) tools

# Context handling
## For multilingual teams

1. In R, start with ellmer and btw

2. When required libraries don't exist in R, try to use reticulate to call specific Python libraries and functions

3*. Create REST APIs for existing R (plumber) or Python (FastAPI) tools

4*. Build a Model Context Protocol (MCP) server

# Context handling
## MCP enables language-agnostic calls to use functions

**Client**

```
Claude Desktop +
Claude Code

Copilot Chat
in VS Code

Positron
Assistant

…
```

**Server**

`mcp_server()`

Tools (R functions), like:

Read R package documentation

Run R code

Inspect objects in the global R environment

(Optionally)

**Session**

`mcp_session()`

Some R Project

Other R Project

A2-Ai

# Context handling
## Calling a Python MCP server in an R session to install an R package



R session to add dplyr

Python MCP server config file

dplyr added and loaded successfully after calling the functions in the MCP server

A2-Ai

EXPLORER

R-AND-AI-TALK-PREP
- Rproj.user
- rv
- .gitignore         U
- .Renviron          U
- .Rhistory
- .Rprofile          U
- llm-calls.qmd      U
- mcp-rv.qmd         U
- r-and-ai-talk-pr... U
- README.md          U
- rproject.toml      U

**mcp-rv.qmd**  U

Preview    ☐ Render on Save    Source    Visual    Insert Code Cell

mcp-rv.qmd > {} (code cell)

```
1   ---
2   title: "mcp"
3   format: html
4   ---
5
    ▷ Run Cell | Run Next Cell
6   ```{r}
7   library(ellmer)
8   library(mcptools)
9   chat <- chat_anthropic(model = "claude-haiku-4-5-20251001")
10  chat$set_tools(mcp_tools())
11  chat$chat(paste("Set working directory to", getwd()))
12  chat$chat("Can you add dplyr and use PPM as the repo if not
    set")
13  ```
14
    ▷ Run Cell | Run Above
15  ```{r}
16  library(dplyr)
17  ```
18
19
```

**rproject.toml**  U

rproject.toml

```
    specifically require a dependency to come from a certain
    repository.
7   # Example: { alias = "PPM", url = "https://packagemanager.
    posit.co/cran/latest" },
8   repositories = [
9
10      |
11
12  # A list of packages to install and any additional
    configuration
13  # Examples:
14      # "dplyr",
15      # {name = "dplyr", repository = "CRAN"},
16      # {name = "dplyr", git = "https://github.com/tidyverse/
        dplyr.git", tag = "v1.1.4"},
17  dependencies = [
18
19      |
20
```

**config.json**  ✕

Users > xufei > .config > mcptools > {} config.json > ...

```
    {
2       "mcpServers": {
3           "rv": {
4               "command": "/Users/xufei/.local/bin/uv",
5               "args": [
6                   "--directory",
7                   "/Users/xufei/project/A2-Ai/mcp/rv",
8                   "run",
9                   "rv.py"
10              ]
11          }
12      }
13  }
14
15      |
```

CONSOLE    TERMINAL    PROBLEMS    OUTPUT    PORTS    DEBUG CONSOLE

📁 ~/project/A2-Ai/r-and-ai-talk-prep

>

# 3 core building blocks of GenAI applications

1. Model Interface
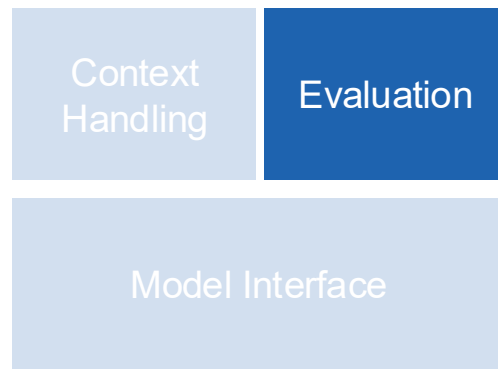
    - Managing LLM API calls

2. Context Handling

    - Handling anything that goes into the LLM

3. Evaluation

    - Validating the results

| Context Handling | Evaluation |
|---|---|
| Model Interface | |

A2-Ai

# Evaluation

Does the application **_systematically_** output desired results?

# Evaluation
## Systematically track inputs, outputs, and correct target responses

> ⌂ > logs > 20250703_145223 > rv > rv_v070_20250630 > 2025-07-03T14-52-24+00-00_rv-v070-20250630_3fxX2tYqYozsX8yNGmYc6D.json     ⓘ

**rv_v070_20250630** anthropic/bedrock/us.anthropic.claude-3-7-sonnet-20250219-v1:0

2025-07-03T14-52-24+00-00_rv-v070-20250630_3fxX2tYqYozsX8yNGmYc6D 📋

|  | ACCURACY | STDERR |
|---|---|---|
|  | 0.846 | 0.067 |

| DATASET | SCORER | CONFIG | DURATION |
|---|---|---|---|
| rv_qna — 13 samples | model_graded_qa | temperature: 0, cache_prompt: true, system: ../prompts/rv_v070_20250630.md, grader: ../prompts/rv_scorer_20250630.md, data_path: ../data/rv_qna.csv, task_name: rv_v070_20250630 | 48.0 sec |

**SAMPLES**   TASK   MODELS   INFO   JSON      FILTER: [_____]   ⓘ SORT: sample asc ▾

| ID | INPUT | TARGET | ANSWER | SCORE |
|---|---|---|---|---|
| 1 | I am migrating to rv - a lot of my packages are from RSPM - is this okay? repositories = [... | **Important:** For A2-AI internal packages, always use PRISM. Your configuration looks good for migration, but you should add PRISM as the last repository in your list:... | **Important:** For A2-AI internal packages, PRISM should be your primary repository source. Your current configuration doesn't include PRISM, which is recommended for better reproducibility and access to intern... | C |
| 2 | renv.lock was migrated to rproject.toml with 2 unresolved packages: scicalc could not be resolved due to:... | Since scicalc and slurmtools are A2-AI internal packages, the version doesn't matter. Simply add them to your dependencies without version specifications:... | **Important:** Since these are A2-AI internal packages, you don't need to worry about the specific version numbers. For internal packages, RV can use the latest available version.... | C |
| 6 | My rv.lock file got edited during a merge of branches. Now I get this error when running rv commands at is is preventing me from... | **Important:** I need to see your lockfile to help properly. This is most likely a result of a merge overwrite that caused the log4r package to be deleted.... | Lockfile Error: Missing Dependencies **Warning: Do not delete your rv.lock file!** This is a common issue that happens during Git merges, and there's a proper way to fix it.... | P |

# Inspect
An open-source framework for large language model evaluations

A2-Ai

# Evaluation – option in R

# Takeaways

1. Today's R GenAI ecosystem is robust enough to do a lot

2. Start building now and refine your requirements as you prototype

3. Bridge strategically if R gaps emerge
   (reticulate, REST APIs, MCP servers)

← model interfaces

← some use cases may require more
   (context handling, evaluation...)

A2-Ai