



s u m m i t

Poland

2024

9-10 OCTOBER 2024, WARSAW

*The Unavoidable AI Revolution
Time for Implementation and Collaboration*

Od dema do produkcyjnych systemów GenAI, czyli o LLMOps

Marek Wiewiórka

GetInData | part of Xebia



Xebia



Marek Wiewiórka

*PhD | Chief Data Architect at
GetInData*



GetInData | Part of Xebia

Experts in **Data, Cloud, Analytics and ML/AI, and GenAI solutions**

Experience in: **media, e-commerce, retail, fintech, banking & telco**

Solution Areas



LLMOps/ML Ops
Platforms



Data, AI & Cloud
Engineering



Stream Processing &
Real-time Analytics



Data Platforms
Modernization &
Migration

Selected technologies



Selected Customers

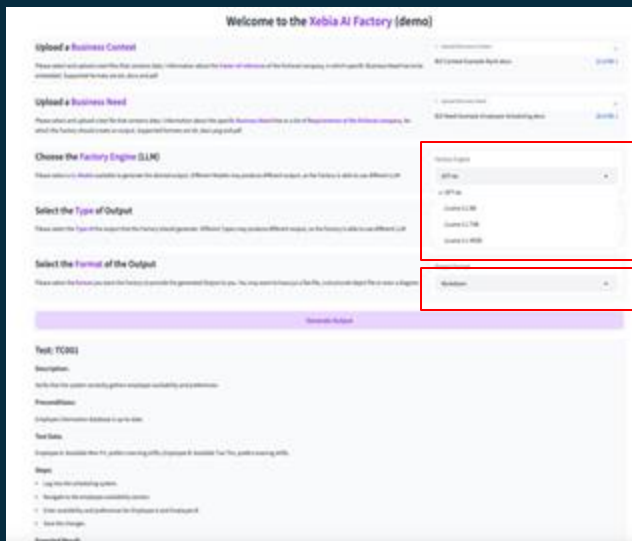


Agenda disclaimer

- Not a comprehensive overview of Large Language Model Operations (*LLMOps*)
- Not (primarily) about tools or platforms
- Not about ready-to-implement processes
- ...but about (a few) **concepts that may help to drive success of GenAI projects**
- **2 different use cases** for inspiration

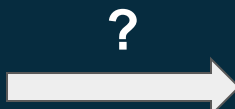
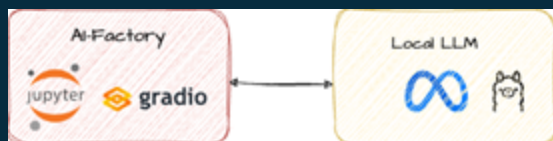
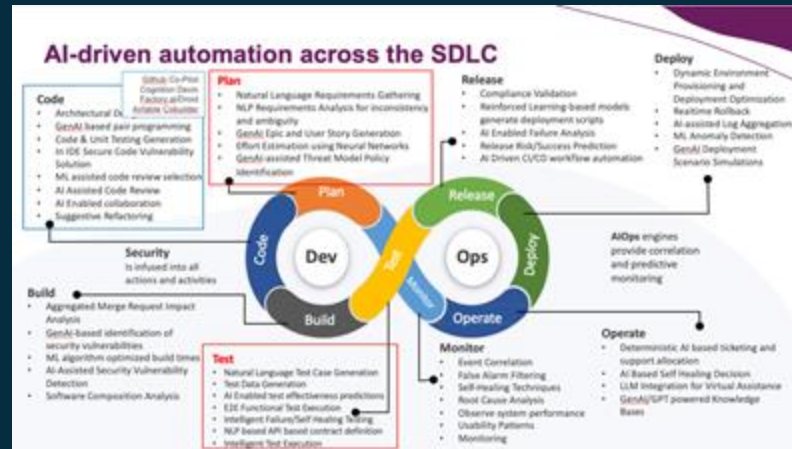


A starting point...



model neutrality

structured output



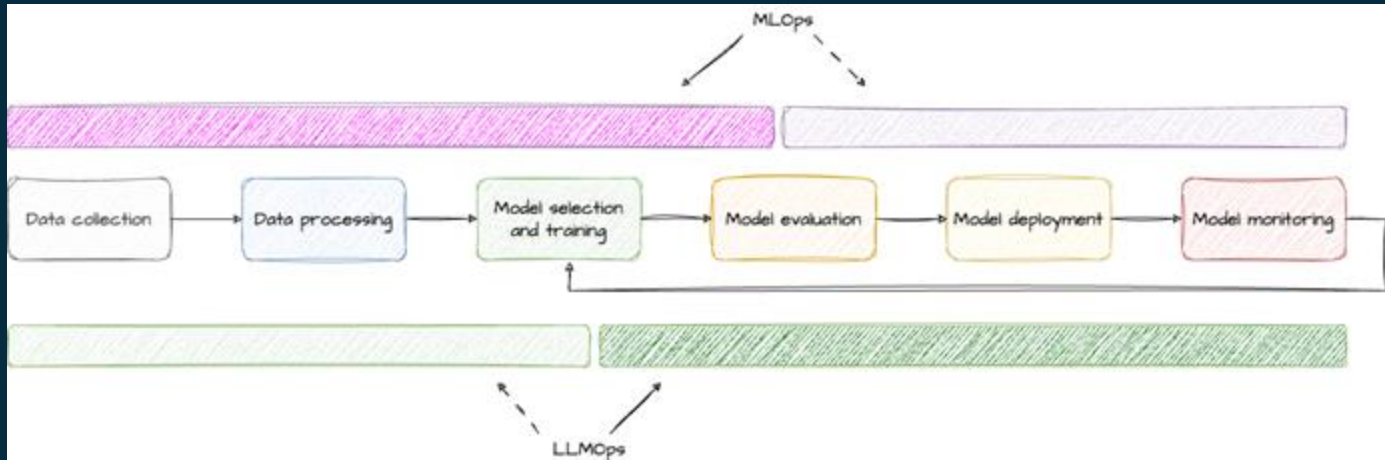
Demo

LLMOps framework!

Production-grade system

MLOps => LLMOps

- Translating ML models into reliable systems and managing their lifecycle
- MLOps and LLMOps have the same goals and principles but differ in **focus**



LLMOps challenges in the enterprise context

- **LLMs churn** - new, more capable models are released frequently
- **Multi model strategies** - optimizing for cost, latency, quality
- **LLM specialization** - one size does not fit all
- On-premise vs managed deployments
- Output of LLM is by nature **non-deterministic**
- Security and data protection
- ... there is not such a thing as LLM backward compatibility out of the box

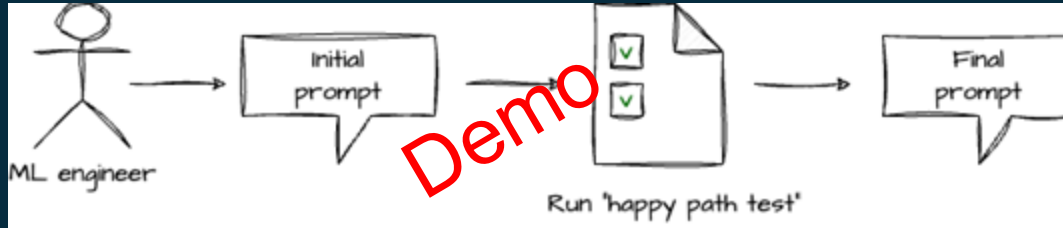
Complexity requires automation

- Prompt optimization
- Controlling LLM Output
- LLM testing and evaluation
- Costs optimization

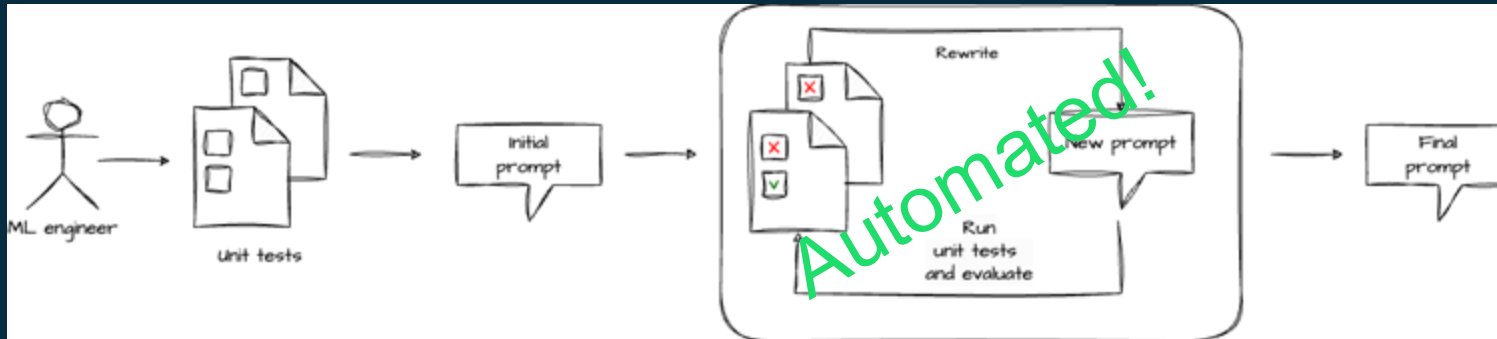
“It’s very easy to make a prototype,” Henley, who studied how copilots are created in his role at Microsoft, says. “It’s very hard to production-ize it.” Prompt engineering—as it exists today—seems like a big part of building a prototype, Henley says, but many other considerations come into play when you’re making a commercial-grade product.

The challenges of making a commercial product include ensuring reliability—for example, failing gracefully when the model goes offline; adapting the model’s output to the appropriate format, because many use cases require outputs other than text; testing to make sure the AI assistant won’t do something harmful in even a small number of cases; and ensuring safety, privacy, and compliance. Testing and compliance are particularly difficult, Henley says, because traditional software-development testing strategies are maladapted for nondeterministic LLMs.

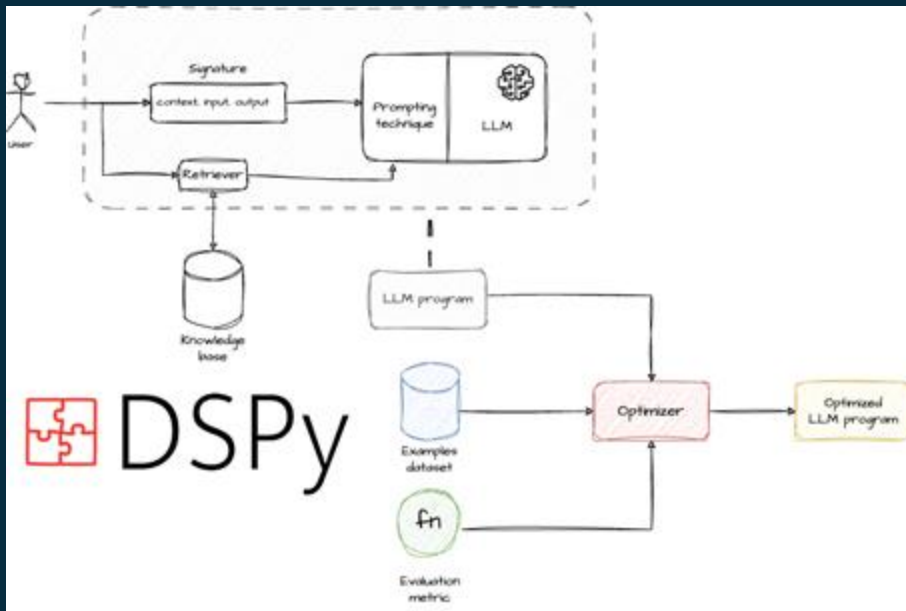
Prompt engineering process



- Similar prompts => different output
- Best prompt specific to a model
- Both instructions and examples (in-context learning) can have great impact on output
- An error-prone and tedious process



Prompt engineering => optimization task



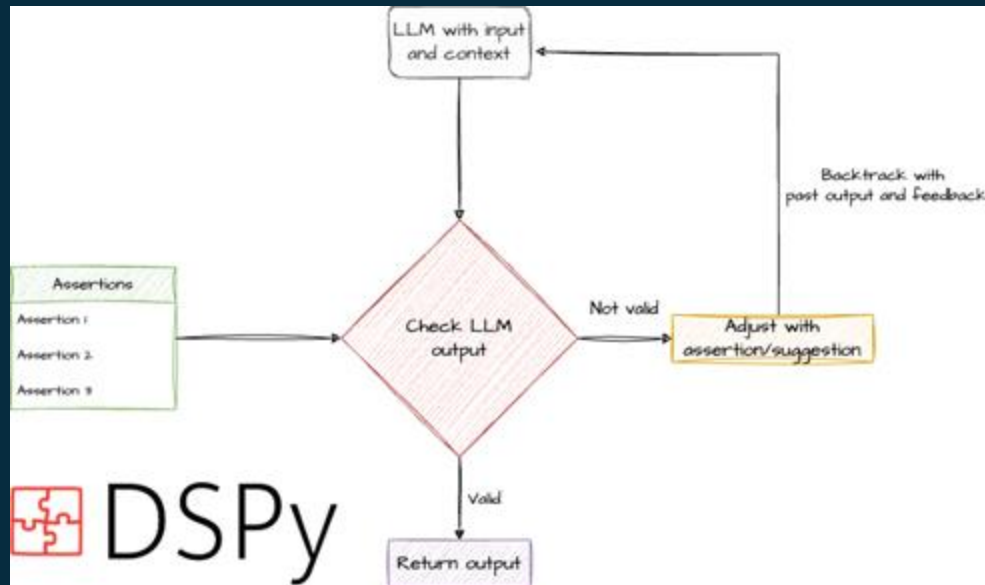
- Different optimizing strategies for both selecting/bootstrapping examples, instructions or models/programs **Ensemble**
- Metric can be an arbitrary function even LLM-based (**LLM-as-a-judge**)
- Can be a **student-teacher** model setup

Zhang, Tuo, Jinyue Yuan, and Salman Avestimehr. "Revisiting OPRO: The Limitations of Small-Scale LLMs as Optimizers." *arXiv preprint arXiv:2405.10276* (2024).

Khatab, Omar, et al. "Dspy: Compiling declarative language model calls into self-improving pipelines." *arXiv preprint arXiv:2310.03714* (2023).

Opsahl-Ong, Krista, et al. "Optimizing Instructions and Demonstrations for Multi-Stage Language Model Programs." *arXiv preprint arXiv:2406.11695* (2024).

Automated strategies for controlling LLM outputs



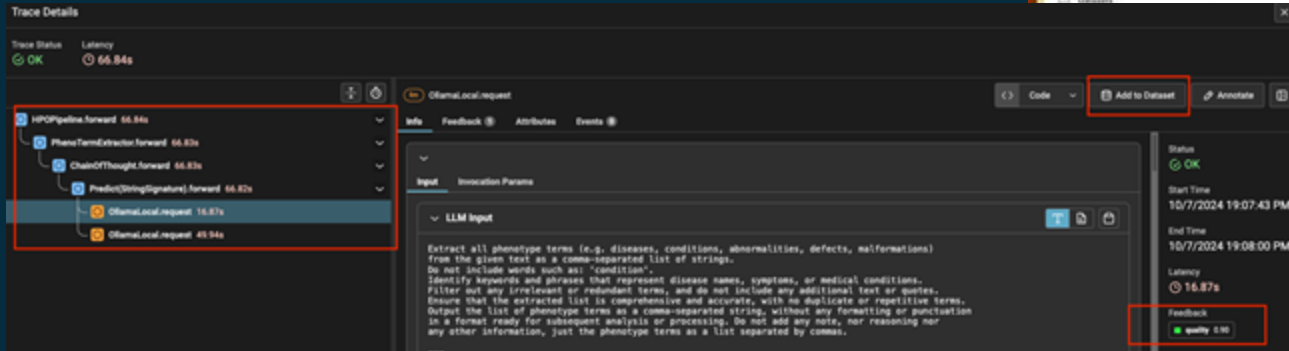
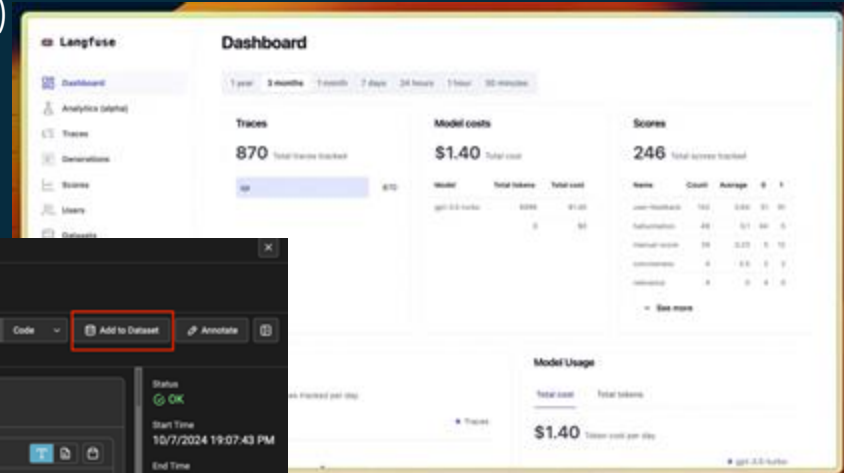
- **Assertions** - general purpose mechanism for guardrailing LLM output
- **Typed Predictions** - specialized for enforcing specific schema using *Pydantic* models for LLM output
- Both can be used in the prompt optimization process

LLM Evaluation

- Absolutely crucial when building a reliable LLM-system
- Depending on the problem can be statistical (e.g. precision, recall, F1) or model-based (LLM-as-a-judge) in more generic cases
- Problem of aligning LLM evaluation with human preferences
 - G-Eval, Prometheus and Evalgen
- Human annotated LLM outputs for calibration
- LLM-assisted criteria and assertion generation

Complexity requires observability

- Open-source tools such as LangFuse, Phoenix and LangSmith emerge, putting high emphasis on LLM observability, including:
 - program metrics, e.g. latency, tokens, costs
 - evaluations scores (optimization process)
 - traces
 - user feedback (annotations)
- user feedback => new datasets



Costs optimization

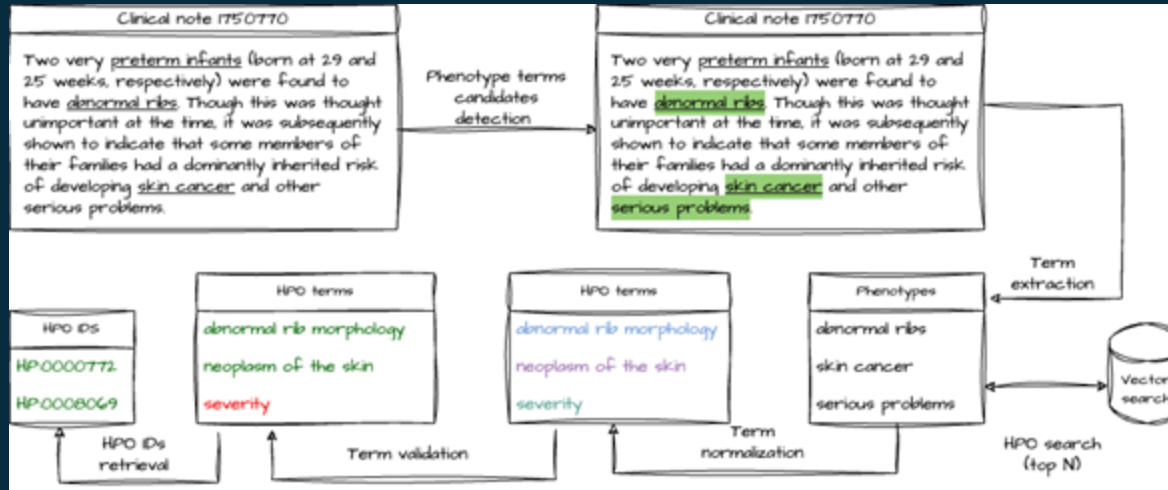
- Different strategies (from most to least complex)
 - hardware optimization
 - LLM optimization (e.g. quantization, scaling down parameters, fine-tuning)
 - LLM routing
 - LLM ensemble optimization, collective wisdom - Mixture-of-Agents
 - prompt optimization
- ... but in order to try to do it you need to have:
 - portable LLM pipelines
 - evaluation datasets and metric functions
 - observability platform

Wang, Junlin, et al. "Mixture-of-Agents Enhances Large Language Model Capabilities." *arXiv preprint arXiv:2406.04692* (2024).

Ong, Isaac, et al. "Routellm: Learning to route llms with preference data." *arXiv preprint arXiv:2406.18665* (2024).

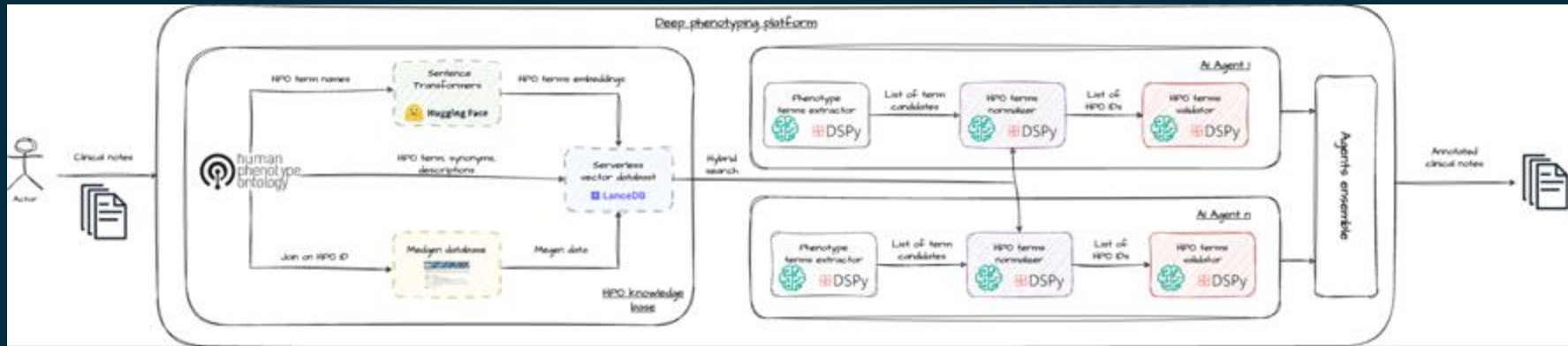
PhenoAgent - use case

- **Deep phenotyping** refers to the comprehensive and detailed analysis of phenotypic traits in organisms
- Two-step procedure involving:
 - concept recognition (finding phenotypic information in the unstructured text) and
 - concept normalization (mapping recognized concepts to the standardized Human Phenotype Ontology (HPO) identifiers)

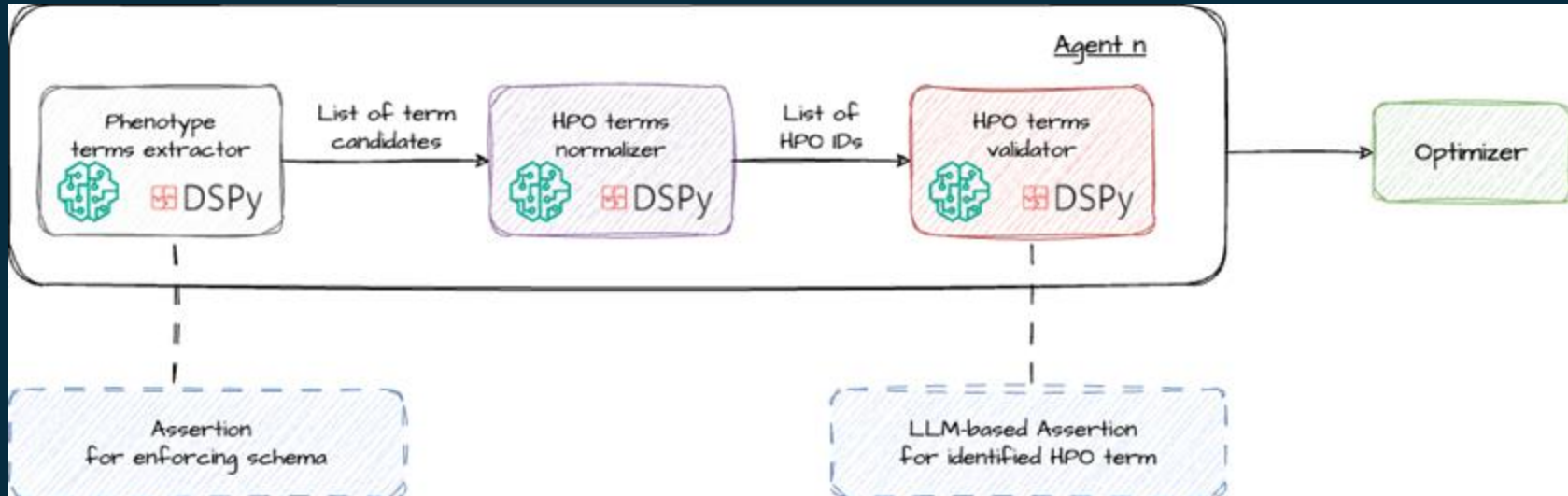


PhenoAgent - architecture

PhenoAgent - first LLM-based tool for an automatic HPO terms annotation powered by RAG and small LLMs ensemble architecture



PhenoAgent - deep dive



PhenoAgent - results and conclusions

- Optimized ensemble of LLM programs of small (and quantized) LLMs can easily **outperform SOTA models** (i.e. Llama-3.1-405/70B)
- RAG architecture can **outperform fine-tuned models** of comparable size
- Using concepts like Assertions and automated prompt optimization help to deliver portable LLM-pipelines
- Using model ensemble and prompt optimization can reduce costs of infrastructure

Tool	Model	Precision	Recall	F1
PhenoGPT	Llama2-7B	0.3136	0.2805	0.2961
PhenoAgent	Llama3-8B	0.5699	0.5511	0.5603
PhenoAgent-MoA-8-3	MoA-8	0.6275	0.6241	0.6258
PhenoAgent-Llama-405	Llama3.1-405B	0.6248	0.5616	0.5915



Xebia

**Want to talk about DATA & AI
with US?**

hello@getindata.com

SCHEDULE A CONSULTATION

