



Motivation

- LLM **agents** have **predefined** sets of **tools**
- Implementing new tools requires **manual work** & technical expertise



Employ agents to **autonomously create new tools** from **research papers** with **code repositories**

Components

Workflow **state** entails **conversation history** and **environment state**

$$s = (h, e) \in \mathcal{H} \times \mathcal{E}$$

There are three types of workflow components (each act on state):

$$\underbrace{\mathcal{H} \times \mathcal{E}}_{\text{old state}} \mapsto \underbrace{\mathcal{H} \times \mathcal{E} \times \mathcal{R}}_{\text{new state return value}}$$

- LLM calls** append to conversation history $\mathcal{H} \mapsto \mathcal{H} \times \mathcal{M}$
- Environment interactions** mutate environment state and return observation $\mathcal{E} \mapsto \mathcal{E} \times \mathcal{O}$ (\mathcal{O} is the set of observations)
- Agents** do both: $\mathcal{H} \times \mathcal{E} \mapsto \mathcal{H} \times \mathcal{E} \times \mathcal{R}$

Problem

Given **paper** + **GitHub repository** + **task description**, generate an LLM-compatible tool



Task

Description: Train a model for biomarker classification using STAMP.
<https://github.com/KatherLab/STAMP>
 (optional full-text article)

Arguments:

- `slide_dir` (str): Path to the folder containing the whole slide images.
 Example: `"/mount/input/TCGA_BRCA"`
- `clini_table` (str): Path to the CSV file containing the clinical data.
 Example: `"/mount/input/clinixlsx"`
- `slide_table` (str): Path to the CSV file containing the slide metadata.
 Example: `"/mount/input/slides.csv"`
- `target_column` (str): Name of the column in `clini_table` that contains the target labels.
 Example: `"pathologic_stage"`

Returns:

- `trained_model` (str): Path to the trained model



Environment definition

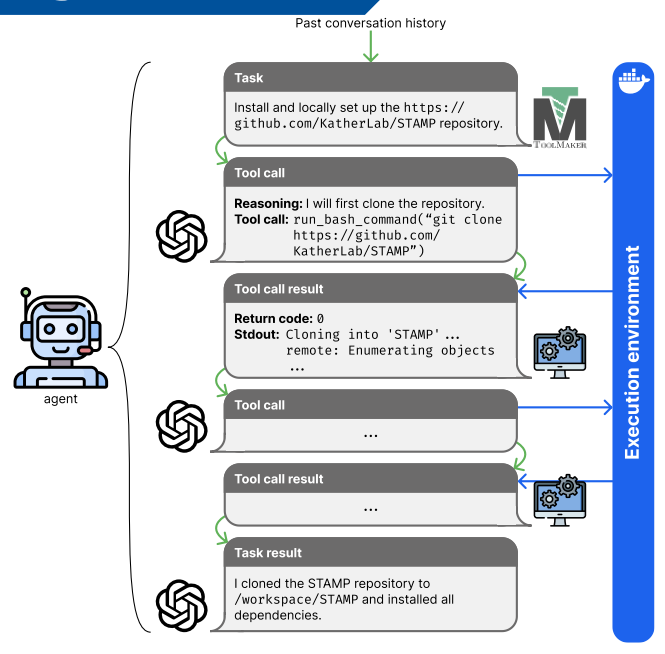
```
FROM python:3.12

RUN git clone https://github.com/KatherLab/STAMP 66 \
  cd STAMP 66 \
  apt update 66 \
  apt install -y openslide-tools 66 \
  pip install -e . 66 \
  stamp init 66 \
  stamp setup
```

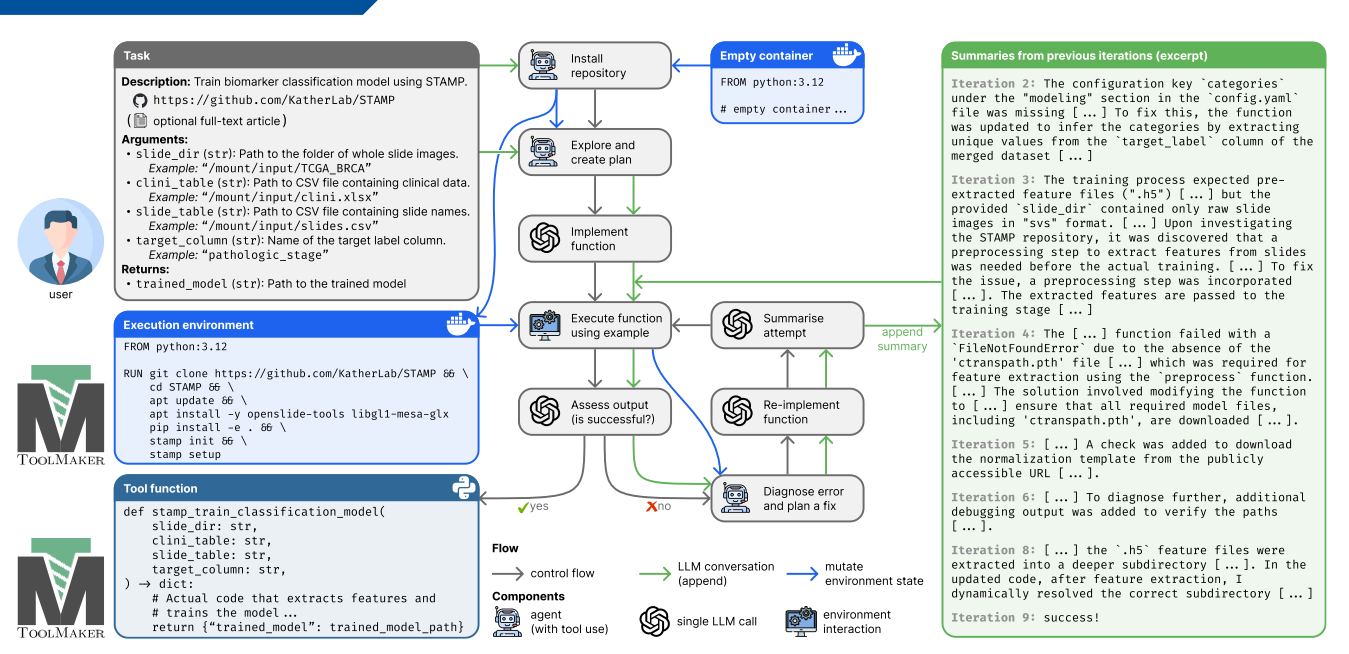
Tool function

```
def stamp_train_classification_model(
    slide_dir: str,
    clini_table: str,
    slide_table: str,
    target_column: str,
) -> dict:
    # Actual code that performs feature
    # extraction and model training...
    return {"trained_model":
            trained_model_path}
```

Agents



Workflow



Main results

Task	TOOLMAKER (ours)					OpenHands (Wang et al., 2024)				
	Invoc.	Tests	Cost	Actions	Tokens	Invoc.	Tests	Cost	Actions	Tokens
Pathology	conch_extract_features (Lu et al., 2024b)	3/3	9/9	\$0.35	15 (10)	171,226	3/3	9/9	\$0.08	5
	musk_extract_features (Xiang et al., 2025)	3/3	6/6	\$1.19	29 (60)	696,386	0/2	4/6	\$0.15	7
	pathfinder_verify_biomarker (Liang et al., 2023)	0/2	4/6	\$0.61	27 (10)	356,825	0/2	4/6	\$0.08	6
	stamp_extract_features (El Nahhas et al., 2024)	3/3	12/12	\$1.12	20 (40)	631,138	0/3	3/12	\$0.07	6
	stamp_train_classification_model (El Nahhas et al., 2024)	3/3	9/9	\$2.27	33 (90)	1,249,521	0/3	0/9	\$0.15	8
Radiology	uni_extract_features (Chen et al., 2024)	3/3	9/9	\$0.61	16 (40)	326,806	0/3	0/9	\$0.25	10
	medsam_inference (Ma et al., 2024)	3/3	6/6	\$0.96	18 (60)	508,954	0/2	0/4	\$0.12	8
Omics	nnunet_train_model (Isensee et al., 2020)	0/2	0/4	\$2.90	35 (90)	1,792,291	0/2	0/4	\$0.12	8
	cytopus_db (Kunes et al., 2023)	3/3	12/12	\$0.41	10 (30)	185,912	0/3	0/6	\$0.36	8
Other	esm_fold_predict (Verkuil et al., 2022; Hie et al., 2022)	2/3	13/15	\$0.66	20 (10)	336,754	0/3	0/6	\$0.11	6
	flowmap_overfit_scene (Smith et al., 2024)	2/2	6/6	\$0.70	18 (50)	358,552	0/3	0/6	\$0.36	15
Other	medsss_generate (Jiang et al., 2025)	3/3	6/6	\$0.53	25 (30)	282,771	3/3	6/6	\$0.15	10
	modernbert_predict_masked (Warner et al., 2024)	3/3	9/9	\$0.66	20 (40)	356,228	0/3	0/6	\$0.13	10
Other	retfound_feature_vector (Zhou et al., 2023)	3/3	6/6	\$0.97	31 (50)	561,936	0/3	0/6	\$0.08	4
	tabpfn_predict (Hollmann et al., 2025)	3/3	9/9	\$0.23	10 (10)	95,257	3/3	9/9	\$0.07	4

Conclusion

Autonomous tool creation is feasible for complex scientific tasks

