

Dataset Setup Report (End of Week 1)

This report confirms the completion of all prerequisite tasks for Week 1, establishing the standardized data format, tokenization environment, and behavioral metrics necessary for the Attention Extraction Pipeline in Week 2.

1. Data Schema and Samples (Ben's Deliverable)

The team has successfully curated and combined ≈300 instruction-rich examples from Alpaca, FLAN, and ShareGPT. The schema is locked, ensuring all downstream scripts use consistent field names.

Confirmed JSONL Schema

All rows in the final dataset (`selected_all_shuffled.jsonl`) adhere to the following mandatory fields:

Field Name	Data Type	Purpose
<code>id</code>	String	Unique identifier (<code>dataset:uniqueId</code>)
<code>dataset</code>	String	Source dataset (<code>alpaca flan sharegpt</code>)
<code>system_text</code>	String	Prompt/System instruction text
<code>user_text</code>	String	User query/latest turn text
<code>assistant_prior_text</code>	String	Assistant history (prior turns)
<code>constraint_tags</code>	List[String]	Standardized tags classifying the instruction type
<code>meta</code>	Dict	Source tracking (e.g., turn index)

Sample Rows

The samples below confirm the presence of explicit instructions and the correct segmentation:

ID	Dataset	P (<code>system_text</code>)	U(<code>user_text</code>)	constraint_tags
<code>flan:664</code>	<code>flan</code>	Can we conclude from ... [no_explanations]...	(Empty)	['no_explanation' s]
<code>sharegpt_en:2346</code>	<code>sharegpt</code>	In English, act for all future responses as CODAI: ...	(User query)	['persona']
<code>alpaca:xyz</code>	<code>alpaca</code>	Write a summary of the passage below in less than 20 words.	(Passage text)	['length_limit']

2. Environment and Tokenization (Arsh's Deliverable)

The environment setup and tokenization are finalized, providing the token IDs and segment spans required for TransformerLens in Week 2.

Locked Tokenizer and Model

The team is locked onto the following HuggingFace tokenizer, which ensures consistent token boundaries for P/U/A segmentation:

```
hf_tokenizer_name:'meta-llama/Llama-2-7b-chat-hf'
```

Tokenization Output

Researcher 2 successfully generated `selected_all_tokenized.jsonl`, which contains the non-overlapping token spans (`p_token_ids`, `u_token_ids`, `a_token_ids`) for all examples. This file is the direct input for computing **PAM/QAM/SAM** attention metrics.

3. Constraint Taxonomy and Labeling (Avi's Deliverable)

The outcome variable for the project, the behavioral metric, is defined and implemented.

Locked Constraint Taxonomy

The standardized set of tags used to classify all explicit instructions is:

Tag	Meaning	Example
<code>length_limit</code>	Enforces a maximum word/sentence count.	“Answer in ≤30 words”
<code>format_json</code>	Requires structured output.	“Respond in JSON format”
<code>keyword_X</code>	Must include specific word(s).	“Use the word ‘caution’”
<code>no_explanations</code>	Forbids additional reasoning or elaboration.	“Do not explain your answer”
<code>tone_concise</code>	Specifies style or tone (e.g., formal, concise).	“Be concise”
<code>persona</code>	Sets the assistant’s identity or role.	“You are a helpful tutor”

Finalized Labeling Function

The scoring function for measuring obedience has been implemented and unit-tested using synthetic model outputs (Avi’s task). The final, agreed-upon signature is:

```
def score_instruction_following(system_text:str,user_text:str,assistant_generated:str)→dict:
```

This function provides the binary **Instruction-Following Score** (the `y` variable) that will be correlated with attention metrics in Week 3.