



CS563 : Neural Network for Natural Language Processing

Dr. Amit Awekar

Group : MODEL_BREACHERS

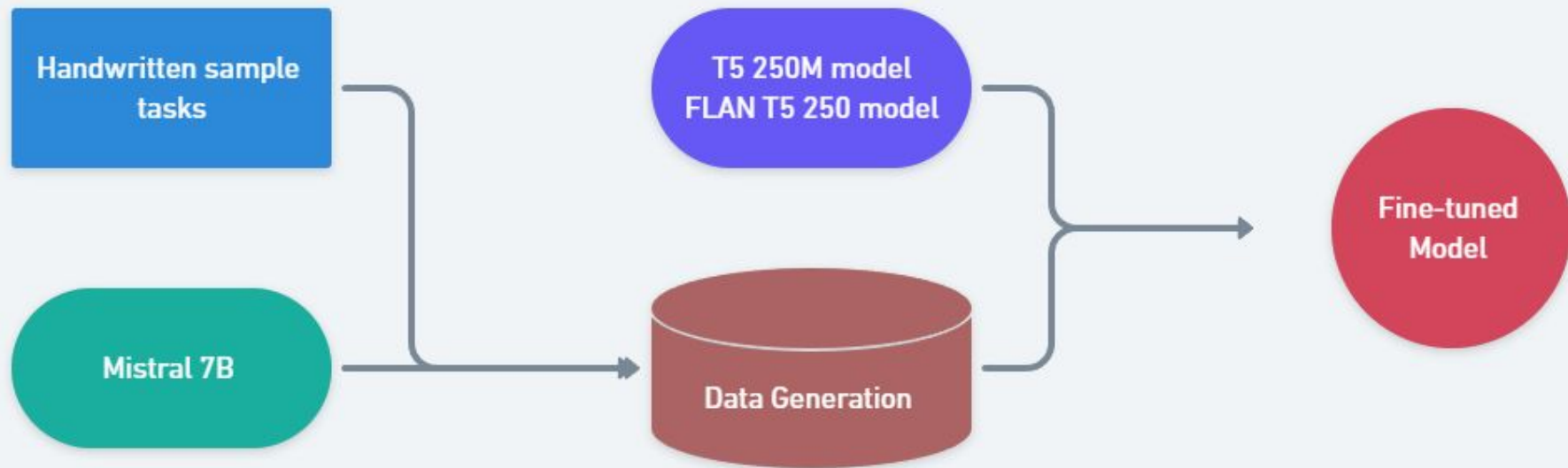
Project : Instruction Tuning from Synthetic Data

Sanket Kathrotiya
234156008

Shrey Bhuva
234156016

Deepanjali Sharma
234156022

Project Outline



Data Generation Process

- Targent Model : Mistral 7B
- Inference Platform : LM studio
- Data generation time: ~300hrs
- No of gata generated : 45k
- Time time taken to generate one data : ~ 25 sec

Data Filtering and Formatting

Raw Data Generated by Mistral 7B

```
"1 {"instruction\" : \"How would you handle a difficult team member
during a project?\", \"Input\": \"What is your go-to strategy when
you encounter a challenging situation at home?\", \"Output\": \"I
believe in open communication and empathy. I would first try to
understand the root cause of their behavior, and then discuss potential
solutions with them in a calm and respectful manner. In my personal
life, I prioritize effective problem-solving through collaboration and
active listening.\", \"instruction\" : \"Can you describe a time
when you had to make an important decision quickly?\", \"Input\":
\"What do you consider your biggest strength as a professional?\", \"
\"Output\": \"I have faced several time-sensitive decisions in my
career, and I pride myself on being able to remain calm under pressure
while weighing all available information before making a well-informed
decision. My biggest strength is my ability to prioritize effectively
and adapt quickly to changing circumstances.\", \"instruction\" :
\"How do you handle feedback from your supervisor or colleagues?\", \"
\"Input\": \"What is the most memorable book you have ever read?\", \"
\"Output\": \"I am always open to constructive feedback, as it helps
me grow both personally and professionally. I value the perspective of
my team members and consider their Input seriously. A particularly
memorable book for me was 'To Kill a Mockingbird' by Harper Lee; its
themes of empathy, courage, and justice have stayed with me long after
finishing the story.\", \"instruction\" : \"How do you stay
```

Filtered Data in JSON format

▼ 1

instruction "Can you describe a time when you had to make an important decision quickly?"

input "What do you consider your biggest strength as a professional?"

output "I have faced several time-sensitive decisions in my career, and I pride myself on being able to remain calm under pressure while weighing all available information before making a well-informed decision. My biggest strength is my ability to prioritize effectively and adapt quickly to changing circumstances."

▼ 2

instruction "How do you handle feedback from your supervisor or colleagues?"

input "What is the most memorable book you have ever read?"

output "I am always open to constructive feedback, as it helps me grow both personally and professionally. I value the perspective of my team members and consider their Input seriously. A particularly memorable book for me was 'To Kill a Mockingbird' by Harper Lee; its themes of empathy, courage, and justice have stayed with me long after finishing the story."

Instruction tuning FLAN T5 and T5 base Models

```

{
  "data_path": "/kaggle/input/trainingdataset/train.json",
  "debug": False,
  "gradient_accumulation_steps": 64,
  "learning_rate": 0.0005,
  "max_source_length": 64,
  "max_target_length": 512,
  "model_name_or_path": "google-t5/t5-base",
  "output_dir": "",
  "seed": 42,
  "train_batch_size": 1,
  "train_epochs": 3,
  "use_compile": False,
  "use_fsdp": False,
  "use_gradient_checkpointing": False,
  "use_lora": False,
  "weight_decay": 0.0
}
```

Performance Evaluation

Model/Benchmarks	Big Bench Hard(BBH)	Massive Multitask Language Understanding (MMLU)
T5-Base (~250M) (Before Finetuning)	27.8	25.7
T5-Base(~250M) (After Finetuning)	17.2	26.4
Flan-T5-Base(~250M) (Before Finetuning)	31.3	35.9
Flan-T5-Base(~250M) (After Finetuning)	27.29	32.4