

FINAL PROJECT FOR
KNOWLEDGE REPRESENTATION AND LEARNING
LLAMA model for NL-FOL task

SKURATIVSKA KATERYNA ID:2081787

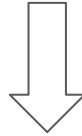
PLAN

1. PROBLEM DEFINITION
2. MODEL SELECTION
3. UNSLOTH + LLAMA
4. LORA
5. MALLS-v0 DATASET
6. DATASET PREPROCESSING
7. MODEL SETUP + LORA SETUP
8. TRAINING SETUP
9. RESULTS
10. CONCLUSIONS

PROBLEM DEFINITION

The aim of the project is to provide a bot model that translates natural-language (NL) statements into first-order logic (FOL) formulas.

If a person is a librarian, they either work in a public library or an academic library.


$$\forall x \text{ (Person}(x) \wedge \text{Librarian}(x) \rightarrow \text{WorkInPublicLibrary}(x) \oplus \text{WorkInAcademicLibrary}(x) \text{)}$$

Model selection

The objective was to provide an effective solution in terms of accuracy, computational complexity, and time efficiency. Therefore, we focused on the LLaMA 8B variant, which offers a balanced trade-off between model complexity and computational efficiency, and integrated it with the Unsloth tool — designed to optimize resource management and streamline the training pipeline — as well as LoRA (Low-Rank Adaptation) into our approach.



with LoRA



unsloth

UNSLOTH LLAMA

Llama model with 8B parameters is less computationally expensive than the 70B parameter model while offering higher accuracy than the lightweight 1-3B versions.

Counting limited computational resources unsloth version of Llama together with Lora technique is a good and efficient solution.

Unsloth Meta-Llama-3.1-8B is optimized versions of Meta's Llama models. The modification was made by applying 4-bit quantization techniques—using tools like the bitsandbytes library — to reduce the memory and computational requirements of these large models. This “trading off” of numerical precision for efficiency does:

- **Lower Memory Footprint:** The 4-bit models require far less GPU memory than their full-precision counterparts.
- **Faster Inference & Training:** Reduced model size enables quicker computations and more cost-effective deployment.
- **Broader Accessibility:** Users can run larger Llama models (even up to 70B parameters) on more modest hardware.

LORA



LoRA (Low-Rank Adaptation) is a fine-tuning technique that improves both training speed and performance by reducing the number of parameters that need to be updated. Here's how it helps:

- **Parameter Efficiency:** Instead of updating all the weights in a large model, LoRA injects small, low-rank trainable matrices into certain layers. This means you only train a fraction of the model's parameters.
- **Faster Training:** With fewer parameters to update, the training process becomes faster and requires less computational power and memory.
- **Better Performance:** By keeping most of the pre-trained model fixed, LoRA reduces the risk of overfitting and preserves the original knowledge, often leading to better generalization on new tasks.

Overall, LoRA makes fine-tuning large models both more resource-efficient and effective.

MALLS-v0 dataset

The MALLS-v0 dataset comprises pairs of natural language statements and their corresponding first-order logic representations. It is provided by the Hugging Face Hub, which allows for easy data uploads via Python and includes a train/test split (27,284 training examples and 1,000 test examples, respectively).

Dataset Viewer	
Auto-converted to Parquet </> API Embed Full Screen Viewer	
Split (2) train · 27.3k rows	
Search this dataset SQL Console	
FOL string · lengths  12 317	NL string · lengths  10 283
$\exists x (Film(x) \wedge ((Drama(x) \wedge LongRuntime(x) \wedge MultipleAwards(x)) \vee (Comedy(x) \wedge ShorterRuntime(x) \wedge ...$	A film can be a drama, have a long runtime, and win multiple awards, or it can be a comedy, have a shorter...
$\forall x (Person(x) \wedge Librarian(x) \rightarrow WorkInPublicLibrary(x) \oplus WorkInAcademicLibrary(x))$	If a person is a librarian, they either work in a public library or an academic library.
$\forall x (HealthySleepHabits(x) \rightarrow ImprovesWellBeing(x))$	Healthy sleep habits improve overall well-being.
$\forall x (Shape(x) \rightarrow (ThreeSides(x) \oplus FourSides(x)))$	A shape can have three or four sides, but not both.
$BoilsAtTemperature(water, 100, seaLevel)$	Water boils at 100 degrees Celsius at sea level.
$\forall x \forall y \forall z (Director(x) \wedge Actor(y) \wedge Screenwriter(z) \wedge Script(z) \rightarrow Collaborates(x, y, z))$	A movie director collaborates with actors and screenwriters to bring a script to life on the screen.
$\forall x \forall y \forall z (OnlineCoursePlatform(x) \wedge (MathematicsSubject(y) \vee LiteratureSubject(y) \vee ...$	An online course platform offers classes in various subjects, such as mathematics, literature, and...
< Previous 1 2 3 ... 273 Next >	

TOKENIZE TEXT

```
malls_prompt = """Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.
```

```
### Instruction:
```

```
{Translate the following natural language statements into their corresponding first-order logic representations.}
```

```
### Input:
```

```
{If a person is a librarian, they either work in a public library or an academic library.}
```

```
### Response:
```

```
{ $\forall x \text{ (Person}(x) \wedge \text{Librarian}(x) \rightarrow \text{WorkInPublicLibrary}(x) \oplus \text{WorkInAcademicLibrary}(x))}$  } """
```


MODEL SETUP + LORA SETUP

```
model, tokenizer = FastLanguageModel.from_pretrained(
    model_name = "unsloth/Meta-Llama-3.1-8B",
    max_seq_length = 2048,
    dtype = None,
    load_in_4bit = True,
)

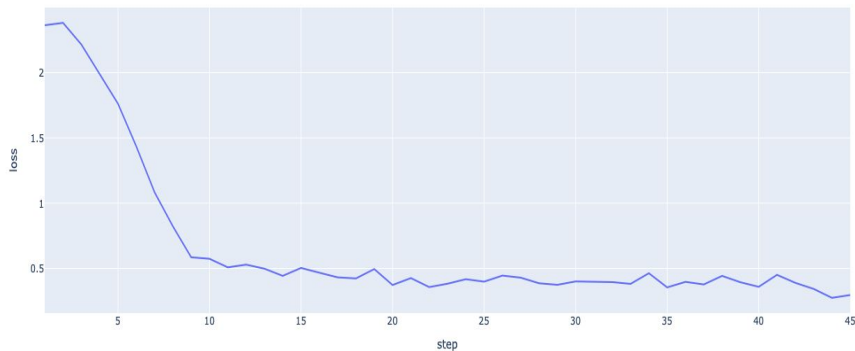
model = FastLanguageModel.get_peft_model(
    model,
    r = 16, # Choose any number > 0 ! Suggested 8, 16, 32, 64, 128
    target_modules = ["q_proj", "k_proj", "v_proj", "o_proj", "gate_proj", "up_proj", "down_proj",],
    lora_alpha = 16,
    lora_dropout = 0, # Supports any, but = 0 is optimized
    bias = "none",    # Supports any, but = "none" is optimized
    # [NEW] "unsloth" uses 30% less VRAM, fits 2x larger batch sizes!
    use_gradient_checkpointing = "unsloth", # True or "unsloth" for very long context
    random_state = 3407,
    use_rslora = False, # We support rank stabilized LoRA
    loftq_config = None, # And LoftQ
)
```

TRAINING SETUP

```
trainer = SFTTrainer(  
    model = model,  
    tokenizer = tokenizer,  
    train_dataset = dataset["train"],  
    dataset_text_field = "text",  
    max_seq_length = max_seq_length,  
    dataset_num_proc = 2,  
    packing = False,  
    args = TrainingArguments(  
        per_device_train_batch_size = 2,  
        gradient_accumulation_steps = 4,  
        warmup_steps = 5,  
        # num_train_epochs = 1,  
        max_steps = 45,  
        learning_rate = 2e-4,  
        fp16 = not is_bfloat16_supported(),  
        bfloat16 = is_bfloat16_supported(),  
        logging_steps = 1,  
        optim = "adamw_8bit",  
        weight_decay = 0.01,  
        lr_scheduler_type = "linear",  
        seed = 3407,  
        output_dir = "outputs",  
        report_to = "none", # Use this for WandB etc  
    ),  
)
```

TRAINING PROCESS

Training Loss Over Time



Before training:

GPU = Tesla T4. Max memory = 14.741 GB.

5.748 GB of memory reserved.

After training:

498.6843 seconds used for training.

8.31 minutes used for training.

Peak reserved memory = 7.01 GB.

Peak reserved memory for training = 1.262 GB.

Peak reserved memory % of max memory = 47.554 %.

Peak reserved memory for training % of max memory = 8.561 %.

CORRECT RESPONSE

A movie is classified as a romantic comedy if it has elements of both romance and comedy.

$\forall x (\text{Movie}(x) \wedge \text{HasRomance}(x) \wedge \text{HasComedy}(x) \rightarrow \text{RomanticComedy}(x))$

$\forall x (\text{Movie}(x) \wedge \text{Romance}(x) \wedge \text{Comedy}(x) \rightarrow \text{RomanticComedy}(x))$

A camera can capture photos or record videos.

$\forall x (\text{Camera}(x) \rightarrow (\text{CapturePhotos}(x) \vee \text{RecordVideos}(x)))$

$\forall x (\text{Camera}(x) \rightarrow (\text{CapturesPhotos}(x) \vee \text{RecordsVideos}(x)))$

A sports event could be played indoors or outdoors, not both.

$\forall x (\text{SportsEvent}(x) \rightarrow (\text{IndoorEvent}(x) \oplus \text{OutdoorEvent}(x)))$

$\forall x (\text{SportsEvent}(x) \rightarrow (\text{PlayedIndoors}(x) \oplus \text{PlayedOutdoors}(x)))$

A garment that is worn on the lower part of the body, has separate openings for each leg, and is made of fabric is a pair of pants or a skirt.

$\forall x (\text{Garment}(x) \wedge \text{WornOnLowerPartOfBody}(x) \wedge \text{SeparateOpeningsForEachLeg}(x) \wedge \text{MadeOfFabric}(x) \rightarrow (\text{Pants}(x) \vee \text{Skirt}(x)))$

$\forall x (\text{Garment}(x) \wedge \text{WornOnLowerBody}(x) \wedge \text{SeparateLegOpenings}(x) \wedge \text{MadeOfFabric}(x) \rightarrow (\text{PairOfPants}(x) \vee \text{Skirt}(x)))$

NOT CORRECT RESPONSE

Chefs prepare meals, waiters serve meals, and customers enjoy meals at a restaurant.

$$\forall x \forall y \forall z (\text{Chef}(x) \wedge \text{Waiter}(y) \wedge \text{Customer}(z) \rightarrow (\text{PreparesMeals}(x) \wedge \text{ServesMeals}(y, z) \wedge \text{EnjoysMeals}(z)))$$
$$\forall x (\text{Chef}(x) \rightarrow \text{PreparesMeals}(x)) \wedge \forall y (\text{Waiter}(y) \rightarrow \text{ServesMeals}(y)) \wedge \forall z (\text{Customer}(z) \rightarrow \text{EnjoysMeals}(z))$$

Winter is cold, and summer is hot.

$$\forall x (\text{Winter}(x) \rightarrow \text{Cold}(x)) \wedge \forall y (\text{Summer}(y) \rightarrow \text{Hot}(y))$$
$$\forall x (\text{Winter}(x) \rightarrow \text{Cold}(x)) \wedge (\text{Summer}(x) \rightarrow \text{Hot}(x))$$

A vacation is relaxing if it includes beautiful scenery and enjoyable activities.

$$\forall x (\text{Vacation}(x) \wedge \text{Relaxing}(x) \rightarrow (\text{BeautifulScenery}(x) \wedge \text{EnjoyableActivities}(x)))$$
$$\forall x (\text{Vacation}(x) \wedge \text{BeautifulScenery}(x) \wedge \text{EnjoyableActivities}(x) \rightarrow \text{Relaxing}(x))$$

An instrument is either a string instrument or a wind instrument, but not both.

$$\forall x (\text{Instrument}(x) \rightarrow ((\text{StringInstrument}(x) \vee \text{WindInstrument}(x)) \wedge \neg(\text{StringInstrument}(x) \wedge \text{WindInstrument}(x))))$$
$$\forall x (\text{Instrument}(x) \rightarrow (\text{StringInstrument}(x) \oplus \text{WindInstrument}(x)))$$

NOT CORRECT RESPONSE

Chefs prepare meals, waiters serve meals, and customers enjoy meals at a restaurant.

$$\forall x \forall y \forall z (\text{Chef}(x) \wedge \text{Waiter}(y) \wedge \text{Customer}(z) \rightarrow (\text{PreparesMeals}(x) \wedge \text{ServesMeals}(y, z) \wedge \text{EnjoysMeals}(z)))$$
$$\forall x (\text{Chef}(x) \rightarrow \text{PreparesMeals}(x)) \wedge \forall y (\text{Waiter}(y) \rightarrow \text{ServesMeals}(y)) \wedge \forall z (\text{Customer}(z) \rightarrow \text{EnjoysMeals}(z))$$

Winter is cold, and summer is hot.

$$\forall x (\text{Winter}(x) \rightarrow \text{Cold}(x)) \wedge \forall y (\text{Summer}(y) \rightarrow \text{Hot}(y))$$
$$\forall x (\text{Winter}(x) \rightarrow \text{Cold}(x)) \wedge (\text{Summer}(x) \rightarrow \text{Hot}(x))$$

A vacation is relaxing if it includes beautiful scenery and enjoyable activities.

$$\forall x (\text{Vacation}(x) \wedge \text{Relaxing}(x) \rightarrow (\text{BeautifulScenery}(x) \wedge \text{EnjoyableActivities}(x)))$$
$$\forall x (\text{Vacation}(x) \wedge \text{BeautifulScenery}(x) \wedge \text{EnjoyableActivities}(x) \rightarrow \text{Relaxing}(x))$$

An instrument is either a string instrument or a wind instrument, but not both.

$$\forall x (\text{Instrument}(x) \rightarrow ((\text{StringInstrument}(x) \vee \text{WindInstrument}(x)) \wedge \neg(\text{StringInstrument}(x) \wedge \text{WindInstrument}(x))))$$
$$\forall x (\text{Instrument}(x) \rightarrow (\text{StringInstrument}(x) \oplus \text{WindInstrument}(x)))$$

CONCLUSION AND FURTHER WORK

Unsloth Llama by itself doesn't have a good evaluation pipeline, but interacting with it helps to check the correctness. However, in the case of FOL, this approach may not be the best solution, because we need not just human-similar response but mathematically correct formula.

Therefore manual evaluation of the equivalence between the FOL ground truth and the model's response can provide a more precise assessment of the results, but it also make it tricky in names-dissimilarities and formulas equivalence.

Another improvement that can be made is train/validation/test split with grid search for the model parameters and training setup. It could be beneficial and provide way better result with more confidence in quality response, although it requires way more computational resources.