

SQuAD 2.0 Question Answering: Experiment Summary

Project Overview

This document tracks the experimental progression of making the Llama-3.2-3B-Instruct model to properly handle unanswerable questions in the SQuAD 2.0 dataset.

Core Challenge: Getting the model to respond with "NO ANSWER" when a question cannot be answered from the provided context, rather than hallucinating an answer.

Test Set: 50 questions from SQuAD 2.0 development set (initial testing), expanding to 1000 questions (final validation), split between answerable (is_impossible = false) and unanswerable (is_impossible = true) questions.

Benchmark Context: According to the SQuAD 2.0 paper:

- State-of-the-art models: 66.3% F1 score
- Human performance: 89.5% F1 score

Project Target: Achieve approximately 89.5% F1 (human-level performance) on the task.

Timeline of Experiments

Phase 1: Establishing Baseline

Initial Baseline Implementation

System Prompt: Minimal concise prompt - "Answer only with text from the context. If the answer is not in the context, reply exactly NO ANSWER".

- Generation Method: Direct model.generate() with deterministic settings.
- Results: 40.0% exact match, 49.3% F1 score on 50 questions.
- Significance: This established the baseline reference point for measuring all subsequent improvements.

Pipeline Infrastructure Testing

- Generation Method: Switched to transformers pipeline for text-generation.
- Key Change: Tested pipeline abstraction which handles message formatting automatically.
- Results: 58.0% exact match, 58.0% F1 score on 50 questions - improved performance over direct generation approach.

Phase 2: Prompt Engineering and Architecture Optimization

Verbose Prompt Exploration

-System Prompt: Enhanced to long, detailed instructions with explicit guidelines about when to answer versus when to say "NO ANSWER".

- Generation Method: Combined verbose prompting with pipeline infrastructure.

- Results: Significant improvement over minimal prompt baseline - balanced performance across both answerable and unanswerable questions.

-Key Finding: Model needed clear, explicit instructions, verbose prompts dramatically outperformed minimal ones.

Phase 3: Dual-Generation Strategy

Self-Consistency via Dual Generation

- Architecture: Implemented dual-generation approach with answer merging logic.

- Primary Generation: Deterministic (do_sample=False) for consistent baseline.

- Secondary Generation**: Stochastic with temperature for diversity.

- Merging Strategy: Only accept answers when both generations agree (normalized comparison).

-Key Innovation: Self-consistency check - if primary and secondary differ, default to "NO ANSWER"

Systematic Configuration Testing

Several systematic experiments were conducted on 50 questions to test different generation strategies:

Config 1: Primary (temperature=0.3) + Secondary (temperature=0.3)

Config 2: Primary (temperature=0.5) + Secondary (temperature=0.5)

Config 3: Primary (deterministic) + Secondary (temperature=0.5)

Key Finding: Configuration testing on 50 questions showed that Config 2 (0.3/0.3) achieved highest F1 on the small sample

Phase 4: GPU Acceleration and Final Optimization

Infrastructure Change: All code in the repository is CPU-based. GPU testing was performed externally on different hardware (**GPU**: NVIDIA GeForce RTX 2080 SUPER (8GB)), dramatically reducing runtime from hours to minutes.

Full Dataset Testing on GPU

- Test Scale: Expanded from 50 to 1000 questions.

- Runtime: ~5-6 minutes for 1000 questions (vs 30-40 minutes for 50 questions on CPU).

- Key Validation: Results closely matched 50-question experiments, confirming the small sample was representative.

Comparative Testing on Full Dataset

Multiple configurations were run on 1000 questions with GPU to compare performance:

-Temperature 0.0 + 0.3

-Temperature 0.3 + 0.5

-Temperature 0.3 + 0.3

-Temperature 0.5 + 0.5

Finding: Configuration 0.3/0.3 showed best balance on the full 1000-question dataset.

Final Enhanced Merging Logic

Prompt Enhancement: Added examples to the system prompt to better illustrate expected behavior.

Merging Logic Refinement :

- If either generation returns "NO ANSWER", final answer is "NO ANSWER".
- Added verification that answer text actually appears in context.

After several additional tests it was concluded that the best temperature configuration is: Primary=0.0 (deterministic), Secondary=0.2 (minimal stochasticity) gave best results

Results on 1000 questions:

(exact', 67.2')

(f1', 70.262')

(total', 1000')

(HasAns_exact', 65.362')

(HasAns_f1', 71.354')

(HasAns_total', 511')

(NoAns_exact', 69.121')

(NoAns_f1', 69.121')

(NoAns_total', 489')

(best_exact', 67.2')

(best_exact_thresh', 0.0')

(best_f1', 70.262')

(best_f1_thresh', 0.0')

time: 367.54 sec

This final configuration achieved the highest F1 score and most balanced HasAns/NoAns performance, demonstrating that conservative answer validation combined with example-driven prompting effectively reduced hallucination while maintaining answer accuracy.

Progress Summary: From baseline of 49.3% F1 to final optimized solution of 70.2% F1 - a 21 percentage point improvement, exceeding the 66.3% F1 benchmark of state-of-the-art models reported in the SQuAD 2.0 paper.

Performance Patterns Observed

What Worked Better

- Verbose prompts with examples outperformed minimal prompts - the model needed clear, explicit instructions and demonstrations.
- Dual-generation with conservative merging helped reduce hallucinations - treating disagreement as uncertainty
- Asymmetric temperature strategy (deterministic primary + low stochastic secondary) balanced consistency with diversity.
- Conservative answer validation - if either generation indicates uncertainty, default to "NO ANSWER".
- Context verification - checking that answer text actually appears in the provided context.
- Batch processing significantly improved throughput over single-question processing.
- GPU acceleration enabled full-dataset testing and rapid iteration.

What Didn't Work

- Symmetric high temperature (0.5/0.5):

The outcome: Too much randomness reduced answer quality.

- Aggressive "Rescue" Strategies:

We attempted to improve recall (the HasAns score) by allowing the Secondary model to provide an answer if the Primary model predicted "NO ANSWER".

The outcome: This approach failed. The Llama-3.2-3B model is prone to hallucinations when "forced" to answer. While this strategy slightly increased the

number of correct answers found, it caused a massive drop in the NoAns accuracy (from ~69% down to ~22% in our tests), as the model began generating plausible sounding but incorrect answers for unanswerable questions.

- Relaxing Context Verification:

We experimented with accepting answers from the model even if they did not appear verbatim in the context passage (relying on semantic similarity).

The outcome: This degraded the Exact Match score. The SQuAD 2.0 metric penalizes paraphrasing.

Explicit Confidence Thresholding:

We hypothesized that asking the model to self-evaluate its certainty could help filter out hallucinations. We modified the system prompt to require a specific output format: [Confidence_Score] Answer Text (e.g., [85] The capital is Paris), and applied a logic filter to discard answers with a score below 70.

The Outcome: This approach caused a significant regression in performance, with the Total F1 score dropping from ~70.5% to ~59.3%. The most dramatic drop was in the NoAns accuracy (which fell to ~40%). The experiment revealed that the Llama-3.2-3B model is highly uncalibrated. It exhibited "confident hallucinations," frequently assigning high confidence scores (90-100) to incorrect answers.

Trade-offs Identified

- Runtime vs Accuracy: Dual-generation nearly doubled processing time but improved confidence in answers.
- HasAns vs NoAns: Some configurations improved unanswerable question handling but slightly reduced answerable question accuracy.

Technical Infrastructure Notes

Execution Environment

- Device :

- All code in repository: CPU-based processing (Intel CPU).
- External GPU testing: Performed on separate hardware with GPU acceleration (NVIDIA GeForce RTX 2080 SUPER (8GB)).
- Precision: bfloat16 for memory efficiency (note: some systems may require float16 /float32 instead depending on hardware support).

- Runtime Comparison:

- CPU: 15-25 minutes for 50 questions (single generation); 40-50 minutes (dual generation).
- GPU (external): 20 - 50 seconds for 50 questions 5-6 minutes for 1000 questions (20x more data in fraction of the time).
- Impact: GPU access enabled full-scale dataset testing and rapid iteration that was not feasible on CPU.

Important Note on Hardware Configuration:

Running time and performance are highly dependent on the specific hardware configuration of the execution environment. It is difficult to predict exact runtime without knowing the processor/GPU capabilities. Configuration parameters should be adjusted based on available hardware:

`'dtype'` (torch.bfloat16, torch.float16, etc.) - depends on hardware support.

`'BATCH_SIZE'` - should be tuned based on available memory.

Different hardware setups may require different configurations for optimal performance.

Data Pipeline

- Input: CSV files from SQuAD 2.0 development set.
- Processing: Streaming row-by-row to maintain memory efficiency.
- Output: CSV with added "final answer" column.
- Evaluation: SQuAD 2.0 metrics (exact match, F1, HasAns/NoAns breakdown).