

**University of Stuttgart** Institute of Industrial Automation and Software Engineering

> Evaluation of Quantized Large Language Models for Semantic Interpretation and Reasoning within Industrial Automation Contexts

> > **Belal Abulabn**

Supervisor: Yuchen Xia



- Motivation
- Basis
- Conception Design
- Implementation
- Evaluation and Verification
- Summary and Outlook



## **Motivation**



#### Rounded Model Size in Millions of Parameters

- 1. High Computational Costs: LLMs require significant computational power for training and inference, making them challenging to deploy in resource-constrained environments.
- 2. Extensive Resource Requirements: Beyond computation, LLMs demand considerable memory and storage, which are scarce on devices with limited processing capabilities.
- **3. High Energy Consumption:** The energy demands of operating LLMs at full scale are substantial

## LLMs and power – Bigger not necessarily better





Billions of bits of data being calculated – moved between memory and processing

Weights, biases and intermediate values moved between memory and processing

More parameters => more data => more bits moved => more energy consumed

## Quantized LLMs for Accessibility and Real-time Processing

- 1. Accessibility and Deployment: Quantization allows LLMs to be more accessible by enabling their deployment on a wider range of devices, including those with limited hardware specifications.
- 2. Energy Efficiency: Reduced resource requirements through quantization lead to lower energy consumption, making LLMs more sustainable and cost-effective for widespread use.
- **3. Real-time Applications**: Quantizing LLMs facilitates real-time processing capabilities on edge devices, crucial for applications requiring immediate responses.







- Motivation
- Basis
- Conception Design
- Implementation
- Evaluation and Verification
- Summary and Outlook



#### Quantization

Quantization is an optimization technique that reduces the precision of the numbers used to represent a model's parameters, which are by default 32-bit floating point numbers. The benefits of this optimization are a smaller model size, better portability, and faster computation.





This method involves converting the weights of a fully trained model to a lower precision format.

QAT integrates the quantization process into the training phase, either during pre-training or fine-tuning





#### Basis







# 8-bit Quantization with LLM.int8()

Optimizing with Precision: Handling Outliers in Quantization Outlier features

- Exceptional features retain their detail through 16-bit floating-point precision (FP16).
- Transition to 8-bit integer (INT8) format streamlines memory usage while upholdin essential precision.





## Adaptive Activation Rounding for Post-Training Quantization



- Simplifies with standard rounding-to-nearest value during post-training quantization.
- Aims to minimize the Frobenius norm, assessing matrix dimensions and calculating error precision.
- Adaptive rounding has been empirically proven to significantly boost model performance and efficiency, particularly in resource-constrained environments such as edge computing and mobil platforms.

- Motivation
- Basis
- Implementation
- Evaluation and Verification
- Summary and Outlook





- To quantize pre-trained models there are different libraries and GitHub repositories that explain how to quantize a model in a specific way.
- Also, there are many quantized models available online which can be used.

#### question, A, B, C, D, correct answer

The wheels and gears of a machine are greased in order to decrease, potential energy, efficiency, output, friction, D, https://github.com/hendrycks/test, arc\_easy.csv, 673, 0.398 xA clean air act must be followed by a, web design course, hard drive manufacturer, math class, computer programming course, B, https://github.com/hendrycks/test, obqa.csv, 1466, 0.385 xWhat is essential for a robot to possess to walk up a flight of stairs?, electricity, skittles, ethics, lava, A, https://github.com/hendrycks/test, obqa.csv, 2764, 0.367 xAs vehicles become more efficient petro consumption, increases, stops, decreases, stay the same, C, https://github.com/hendrycks/test, obqa.csv, 3511, 0.343 You, 2 months ago • Un

- Objective: To assess the effectiveness of quantized language models in industrial automation applications.
- Dataset Foundation: Utilized a curated set of 875 single-choice questions, extracted from a broader dataset aimed at "MEASURING MASSIVE MULTITASK LANGUAGE UNDERSTANDING."
- Selection Criteria: Questions were carefully chosen to specifically evaluate the models' comprehension and problem-solving abilities within the industrial automation domain.
- Purpose: This specialized dataset provides a critical benchmarking tool for measuring the nuanced understanding and reasoning capabilities of quantized models in a targeted industrial context.

#### Implementation

#### question,A,B,C,D,correct answer

The wheels and gears of a machine are greased in order to decrease, potential energy, efficiency, output, friction, D, https://github.com/hendrycks/test, arc\_easy.csv, 673, 0.398 xA clean air act must be followed by a, web design course, hard drive manufacturer, math class, computer programming course, B, https://github.com/hendrycks/test, obqa.csv, 1466, 0.385 xWhat is essential for a robot to possess to walk up a flight of stairs?, electricity, skittles, ethics, lava, A, https://github.com/hendrycks/test, obqa.csv, 2764, 0.367 xAs vehicles become more efficient petro consumption, increases, stops, decreases, stay the same, C, https://github.com/hendrycks/test, obqa.csv, 3511, 0.343 You, 2 months ago \* Unc



- Motivation
- Basis
- Implementation
- Evaluation and Verification
- Summary and Outlook







## Llama-2-7b-chat and Its Quantized Models (q2k, q3kl, q5kl, q8)



Correct Answers Wrong Answers

## LLM.int8() vs q8



- Ensured fair comparison with consistent parameters such as temperature and prompts.
- Hypothesized q8's potential for higher performance due to alignment with float16 and computational strategy.
- q8 outperforms int8 in accuracy, resulting in fewer invalid responses, confirming
- improved model understanding.

## **Summary and Outlook**

- Investigated quantization on LLMs, emphasizing semantic interpretation for industrial automation.
- Utilized Llama-2-7b-chat model to evaluate efficiency in technical reasoning within industrial contexts.
- Benchmarked against 857 questions from industrial automation papers, measuring accuracy against model size reduction.
- Established that quantization effectively balances model size with performance, maintaining substantial accuracy.

- Confirmed quantized LLMs' viability in industrial automation, interpreting complex data accurately.
- q8 quantization method notably outperformed int8, underscoring the importance of strategic quantization selection.



**University of Stuttgart** Institut of Industrial Automation and Software Engineering

## Thank you!



#### **Belal Abulabn**

e-mail st18211@stud.uni-stuttgart.de phone +49 (0) 711 685fax +49 (0) 711 685-

University of Stuttgart

