

Comparative Data on Llama 3 Quantization Methods

Author: Smolecule Technical Support Team. Date: February 2026

Compound Focus: IQ-3

Cat. No.: S1944062

Get Quote

The following table summarizes the performance of different quantization methods for the **Llama 3 8B** model on the MMLU (Massive Multitask Language Understanding) test, which is a common benchmark for evaluating AI model capabilities [1].

Model Size (GB)	MMLU Score (%)	Bits per Weight (bpw)	Quantization Method	Framework
13.98	65.20	16.00	FP16 (baseline)	GGUF / ExL2
6.99	64.53	8.00	8-bit	Transformers
5.73	65.06	6.56	Q6_K	GGUF
5.00	64.90	5.67	Q5_K_M	GGUF
4.30	64.64	4.82	Q4_K_M	GGUF
3.87	64.39	4.28	IQ4_XS	GGUF
3.53	62.89	3.79	Q3_K_M	GGUF
3.49	63.42	4.00	4-bit NF4	Transformers
3.31	62.55	3.50	IQ3_M	GGUF

Model Size (GB)	MMLU Score (%)	Bits per Weight (bpw)	Quantization Method	Framework
3.23	60.28	3.50	IQ3_XS	GGUF

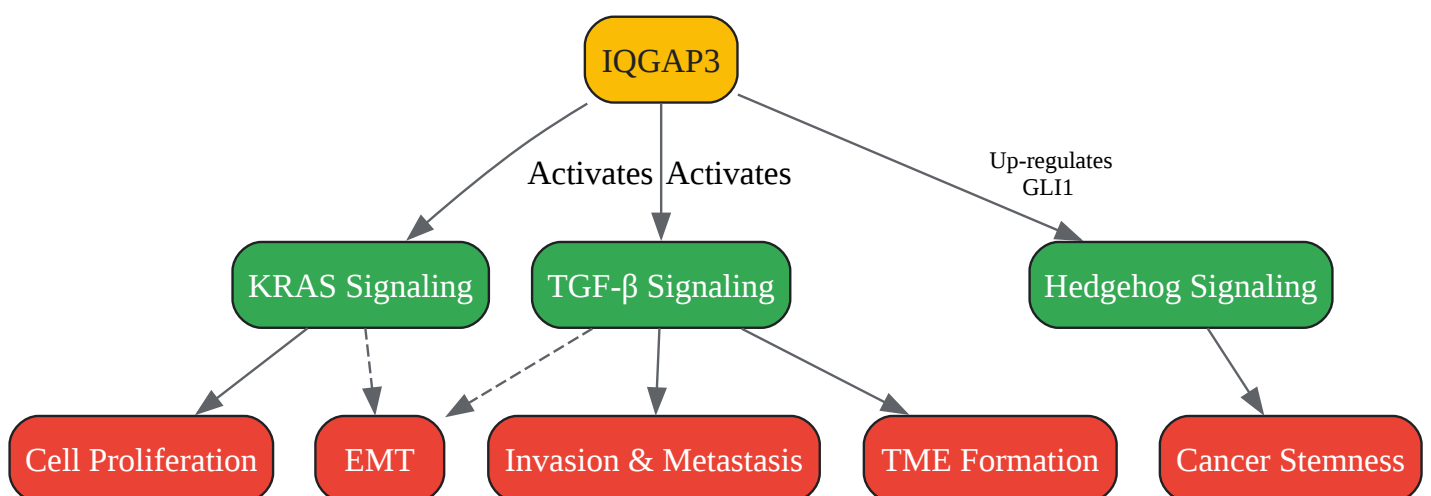
Key Takeaways from the Data:

- **Performance Trade-off:** As the model is compressed more heavily (lower bpw and smaller file size), the MMLU score generally decreases, demonstrating a trade-off between efficiency and performance [1].
- **IQ3_M Performance:** The IQ3_M method offers a balance, reducing the model size to 3.31 GB while maintaining 62.55% accuracy on the MMLU, which is competitive with other 3-4 bit methods [1].

Understanding IQGAP3 in Cellular Signaling

In biomedical research, **IQGAP3** is not a method but a scaffolding protein that is overexpressed in various cancers and plays a crucial role in regulating multiple signaling pathways that drive tumor growth and metastasis [2] [3] [4].

The diagram below synthesizes findings from recent studies to show how IQGAP3 acts as a central hub in a network that promotes cancer malignancy.



[Click to download full resolution via product page](#)

Experimental Insights into IQGAP3 Function:

Research into IQGAP3 employs standardized molecular biology techniques. Key experimental approaches and findings include [2] [3]:

- **Gene Knockdown:** Using small interfering RNA (siRNA) to inhibit IQGAP3 expression in cancer cell lines (e.g., AGS, NUGC3, A549, H1299) is a primary method to study its function.
- **Transcriptomic Analysis:** RNA sequencing (RNA-seq) following IQGAP3 knockdown reveals downregulation of key pathways, including **KRAS signaling** and **TGF- β signaling**.
- **Phenotypic Assays:** Functional experiments show that IQGAP3 depletion reduces cancer cell **proliferation, migration, invasion, and spheroid colony formation** (a proxy for cancer stemness).
- **Protein Interaction:** Co-immunoprecipitation (Co-IP) experiments have identified physical interactions between IQGAP3 and other proteins, such as **RAD17** in lung cancer and **GLI1** in the Hedgehog pathway.
- **In Vivo Validation:** Xenograft models in immunodeficient mice demonstrate that IQGAP3 knockdown suppresses **tumor growth and lung metastasis**.

How to Proceed with Your Comparison Guide

Given the distinct nature of the two "IQ-3" subjects, here is some guidance for your project:

- **For a comparison of AI model quantization:** The data for IQ3_M and other methods is well-suited for a technical guide. You can expand the comparison by including metrics beyond MMLU, such as inference speed and resource usage on different hardware.
- **For a comparison in a drug development context:** A "cross-method comparison" for **IQGAP3** would involve evaluating different techniques to target this protein (e.g., small molecule inhibitors, siRNA, antibody-based therapies). The current literature strongly supports its role as a promising **therapeutic target** across multiple cancers, including gastric, lung, and pancreatic cancer [2] [3] [4].

Need Custom Synthesis?

Email: info@smolecule.com or [Request Quote Online](#).

References

1. GitHub - matt-c1/llama- 3 -quant- comparison : Comparison of the... [github.com]

2. IQGAP3 signalling mediates intratumoral functional ... [pmc.ncbi.nlm.nih.gov]
3. IQGAP3 activates Hedgehog signaling to confer stemness ... [nature.com]
4. IQ Motif Containing GTPase-Activating Protein 3 Is Associated ... [pmc.ncbi.nlm.nih.gov]

To cite this document: Smolecule. [Comparative Data on Llama 3 Quantization Methods]. Smolecule, [2026]. [Online PDF]. Available at: [<https://www.smolecule.com/products/b1944062#iq-3-cross-method-comparison>]

Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While Smolecule strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

Technical Support: The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [Contact our Ph.D. Support Team for a compatibility check]

Need Industrial/Bulk Grade? Request Custom Synthesis Quote

Smolecule

Your Ultimate Destination for Small-Molecule (aka. smolecule) Compounds, Empowering Innovative Research Solutions Beyond Boundaries.

Contact

Address: Ontario, CA 91761, United States
Phone: (512) 262-9938
Email: info@smolecule.com
Web: www.smolecule.com