



**ICLR**

# The Geometry of LLM Quantization: GPTQ as Babai's Nearest Plane Algorithm

Jiale Chen ([jiale.chen@ist.ac.at](mailto:jiale.chen@ist.ac.at))

Yalda Shabanzadeh

Elvir Crnčević

Torsten Hoefler

Dan Alistarh



Institute of  
Science and  
Technology  
Austria



**Red Hat**

**ETH** zürich

# Background & Motivation

- LLMs are expensive to store and run
- Post-training quantization reduces weights from 16-bit to low bitwidth
- GPTQ is widely used, but its success was not well understood theoretically
- **Main contribution:** GPTQ can be interpreted geometrically as Babai's nearest plane algorithm for the closest vector problem (CVP) on a lattice

# From LLM Quantization to Lattice Problem

## LLM Linear Layer Quantization

Find  $Q$  or  $\mathbf{z}_i$  to minimize

$$\|\mathbf{X}Q - \mathbf{X}W\|_F^2 = \sum_i \|\mathbf{X} \text{diag}(\mathbf{s}_i) \mathbf{z}_i - \mathbf{X}w_i\|^2$$

## Closest Vector Problem (CVP)

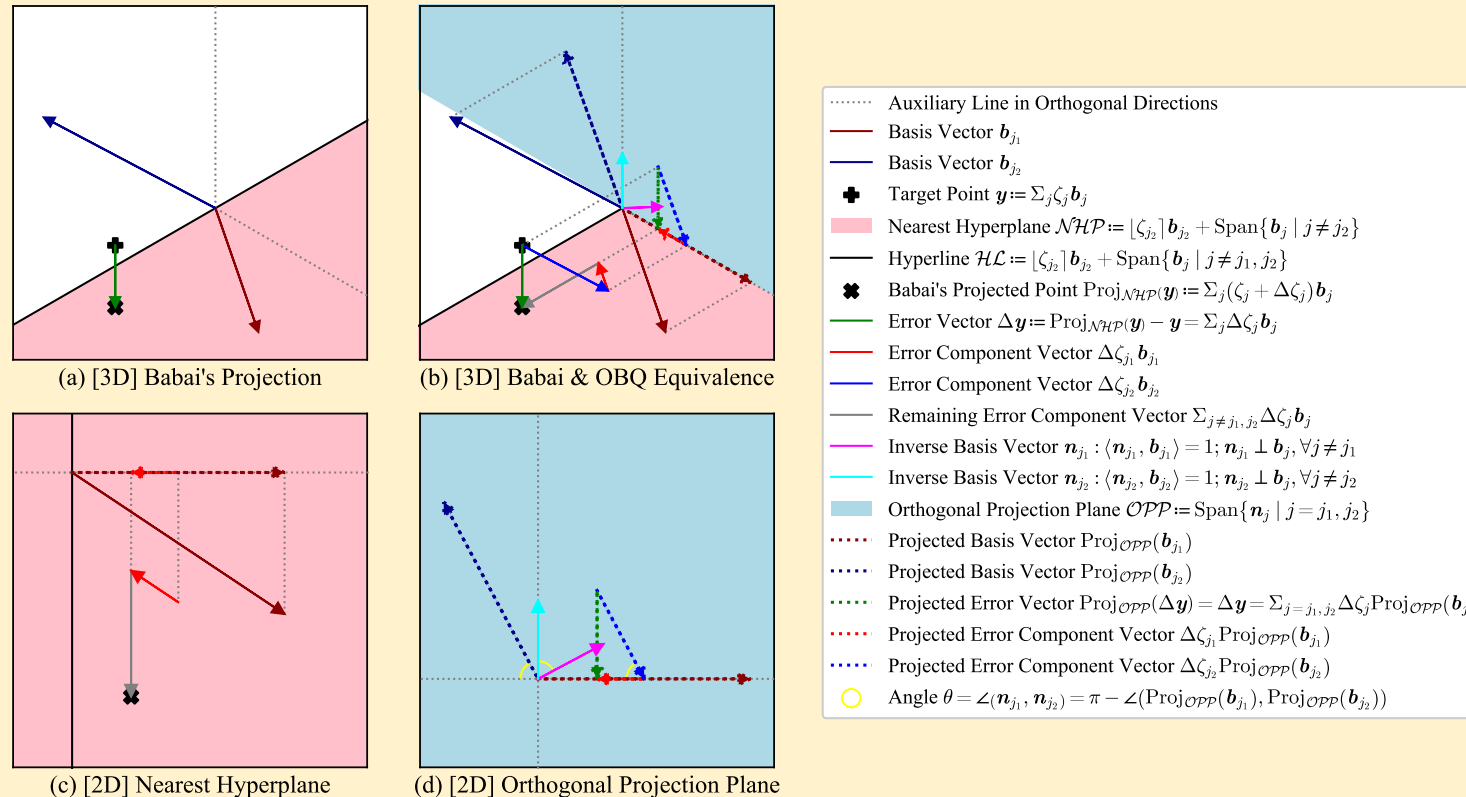
Find  $\mathbf{z}$  to minimize

$$\|\mathbf{B}\mathbf{z} - \mathbf{y}\|^2$$

- Activations  $\mathbf{X}$  → Basis directions
- Scales  $\mathbf{s}_i$  → Basis stretches
- Quantized integers  $\mathbf{z}_i$  → Lattice coordinates  $\mathbf{z}$
- Target  $\mathbf{X}w_i$  → Point to approximate  $\mathbf{y}$

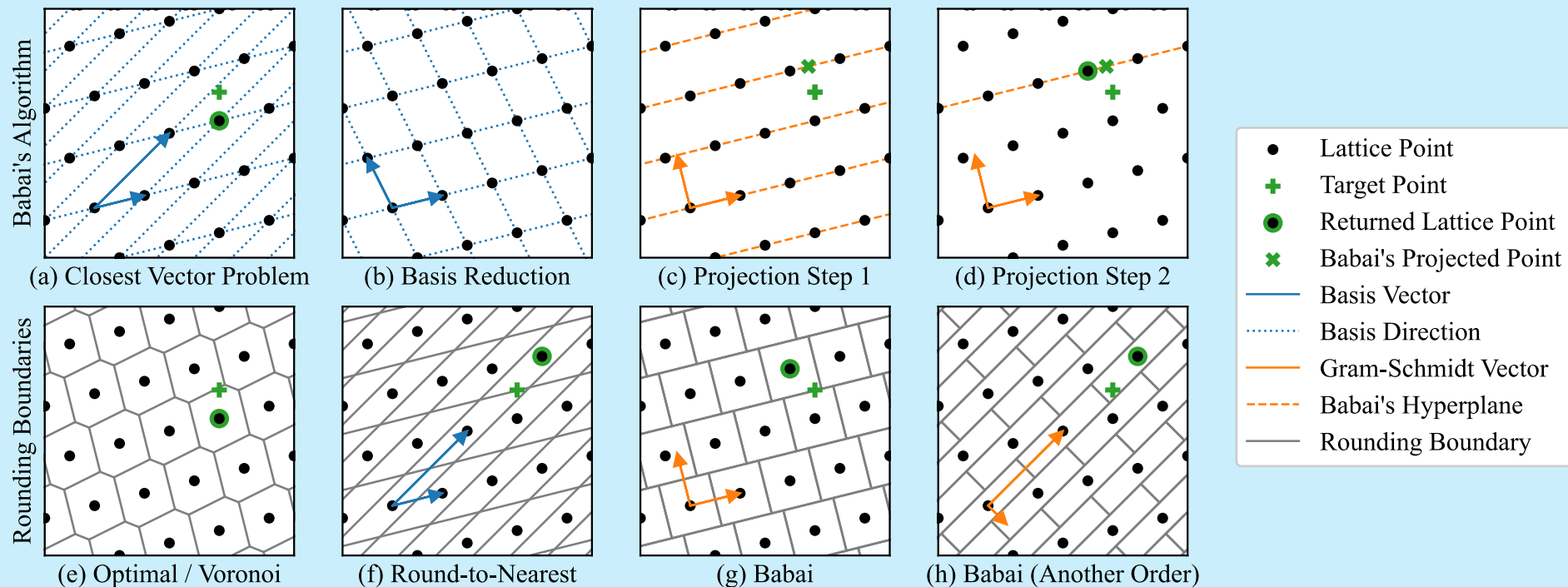
# Main Theorem: GPTQ is Babai's Algorithm

- OBQ/GPTQ error propagation  $\Leftrightarrow$  Projection onto the nearest hyperplane
- Geometrically, GPTQ performs an orthogonal walk through affine subspaces

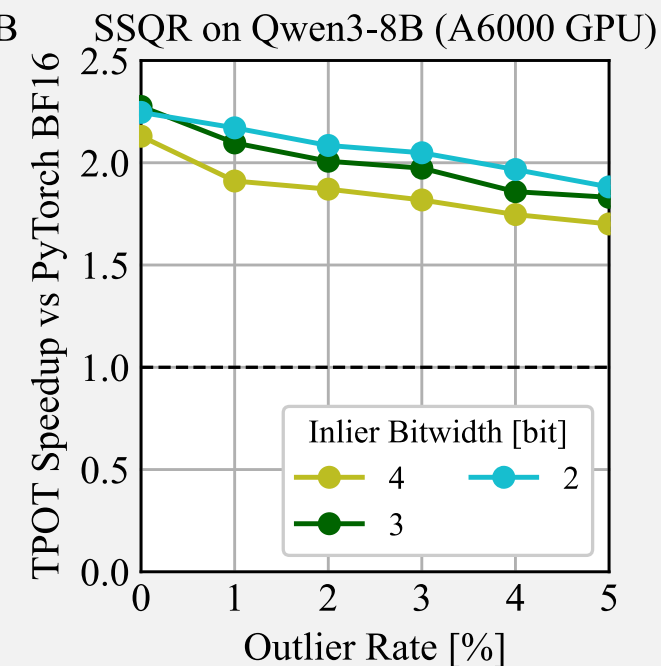
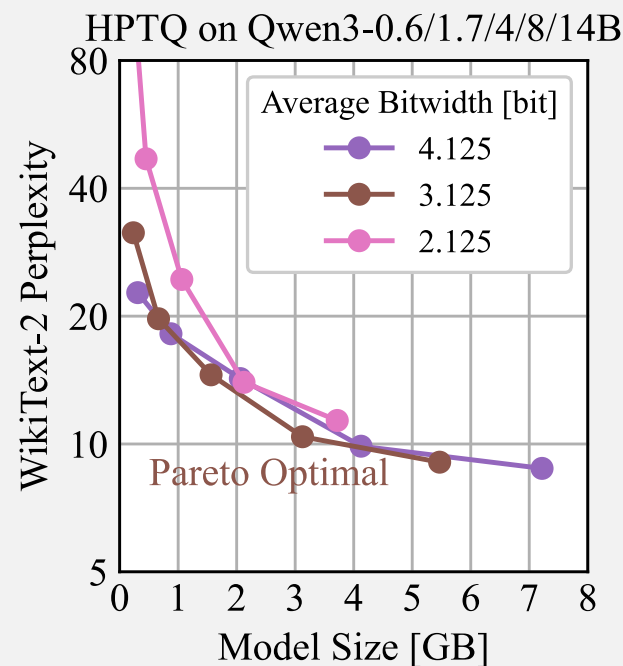
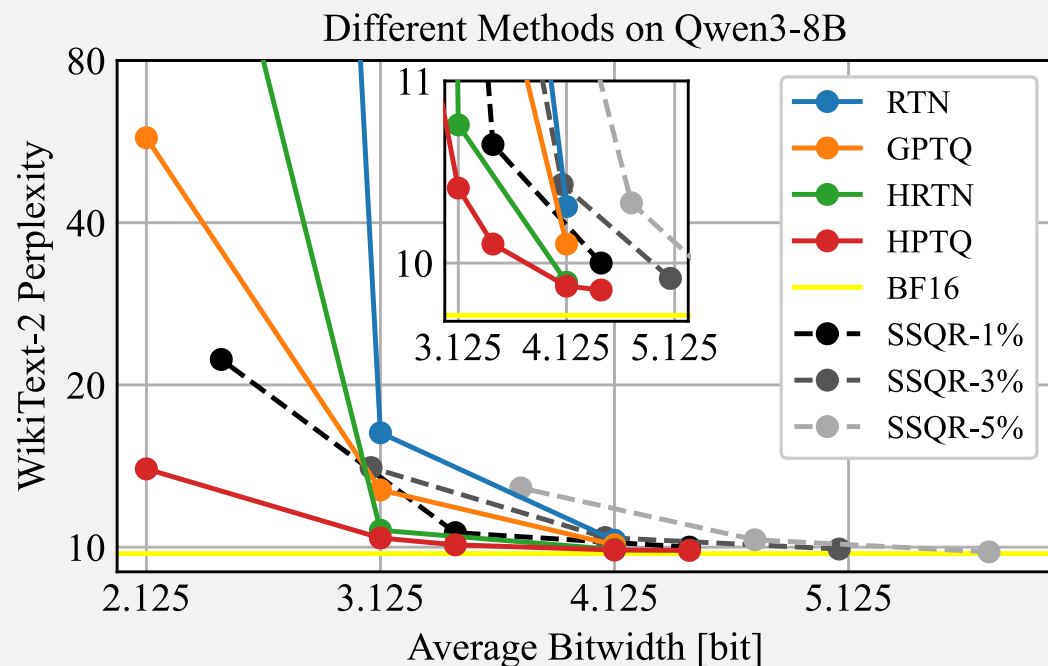


# What the Equivalence Brings: Error Bounds

- In the no-clipping setting, GPTQ inherits Babai's tight error bound
- The bound depends on the LDL decomposition of the permuted Hessian
- Quantization order matters



# Practical Impact: Better Quantizations



Paper  
arXiv:2507.18553



Code  
IST-DASLab/GPTQ-Babai