



Paper

arXiv:2507.18553 IST-DASLab/GPTQ-Babai



Code

# The Geometry of LLM Quantization:

## GPTQ as Babai's Nearest Plane Algorithm

Jiale Chen (jiale.chen@ist.ac.at),

Yalda Shabanzadeh, Elvir Crnčević,

Torsten Hoefler, Dan Alistarh

### Equivalence

#### LLM Linear Layer Quantization

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

$$\|XQ - XW\|_F^2 = \sum_i \|\mathbf{X} \text{diag}(\mathbf{s}_i) \mathbf{z}_i - \mathbf{X} \mathbf{w}_i\|^2$$



#### Closest Vector Problem (CVP)

Find  $\mathbf{z}$  to minimize

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

#### Algorithm 1: GPTQ

Input:  $W, S, X, P, \lambda, \mathbb{Z}_\dagger$   
Output:  $Z, Q$

- $H \leftarrow P^\top (X^\top X + \lambda I) P$
- $L \leftarrow \text{LDL}(H^{-1})$
- $W, S \leftarrow P^{-1}W, P^{-1}S$
- $Q, Z \leftarrow W, 0$
- for  $j \leftarrow 1$  to  $c$  do
  - $\zeta \leftarrow W[j, :]/S[j, :]$
  - $Z[j, :] \leftarrow \text{ROUND}(\zeta, \mathbb{Z}_\dagger)$
  - $Q[j, :] \leftarrow Z[j, :] * S[j, :]$
  - $\epsilon \leftarrow Q[j, :] - W[j, :]$
  - $W[j :, :] \leftarrow W[j :, :] + L[j :, j]\epsilon$
- end
- $Z, Q \leftarrow PZ, PQ$

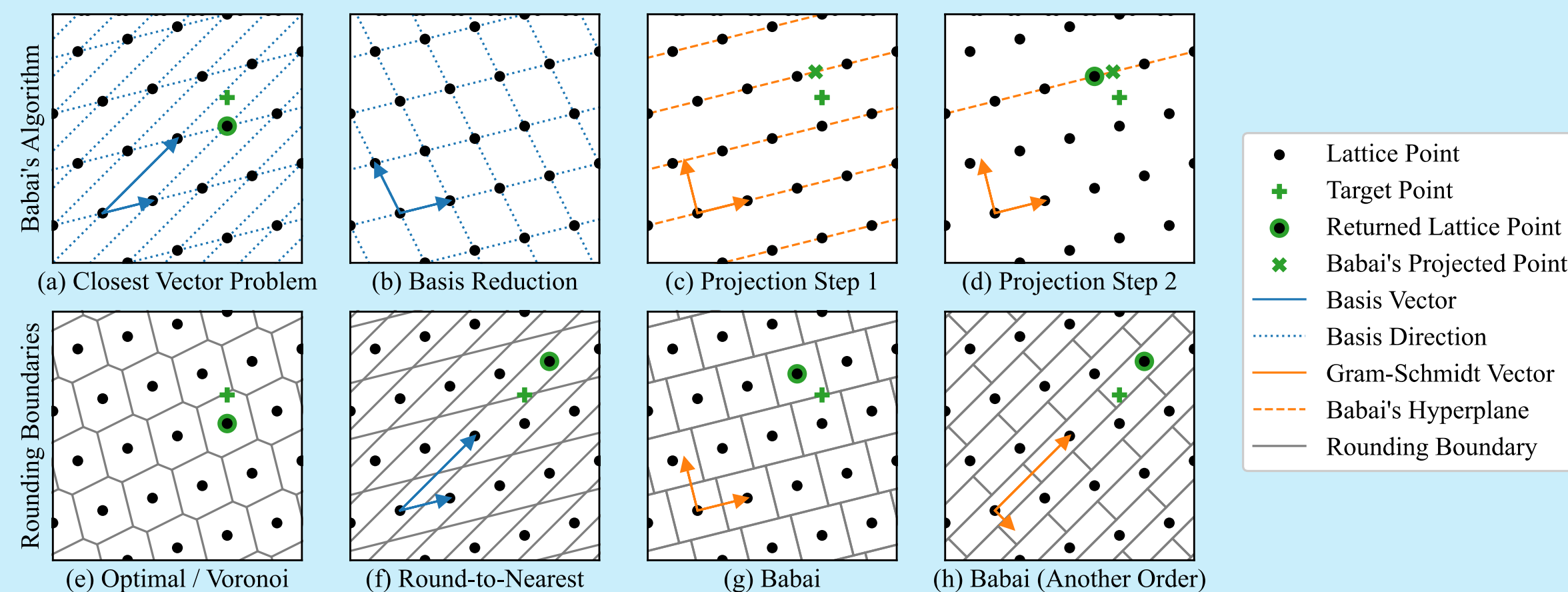
#### Algorithm 2: Babai's Nearest Plane

Input:  $B, y$   
Output:  $z$

- $T \leftarrow \text{LLL}(B)$  // transformation
- $A \leftarrow BT$  // basis reduction
- $\Phi \leftarrow \text{QR}(A)$  // orthogonalize
- $y', z \leftarrow y, 0$
- for  $j \leftarrow c$  to 1 do
  - $\zeta \leftarrow \langle \Phi[:, j], y' \rangle / \langle \Phi[:, j], A[:, j] \rangle$
  - $z[j] \leftarrow \text{ROUND}(\zeta, \mathbb{Z})$
  - $y' \leftarrow y' - A[:, j]z[j]$
- end
- $z \leftarrow Tz$

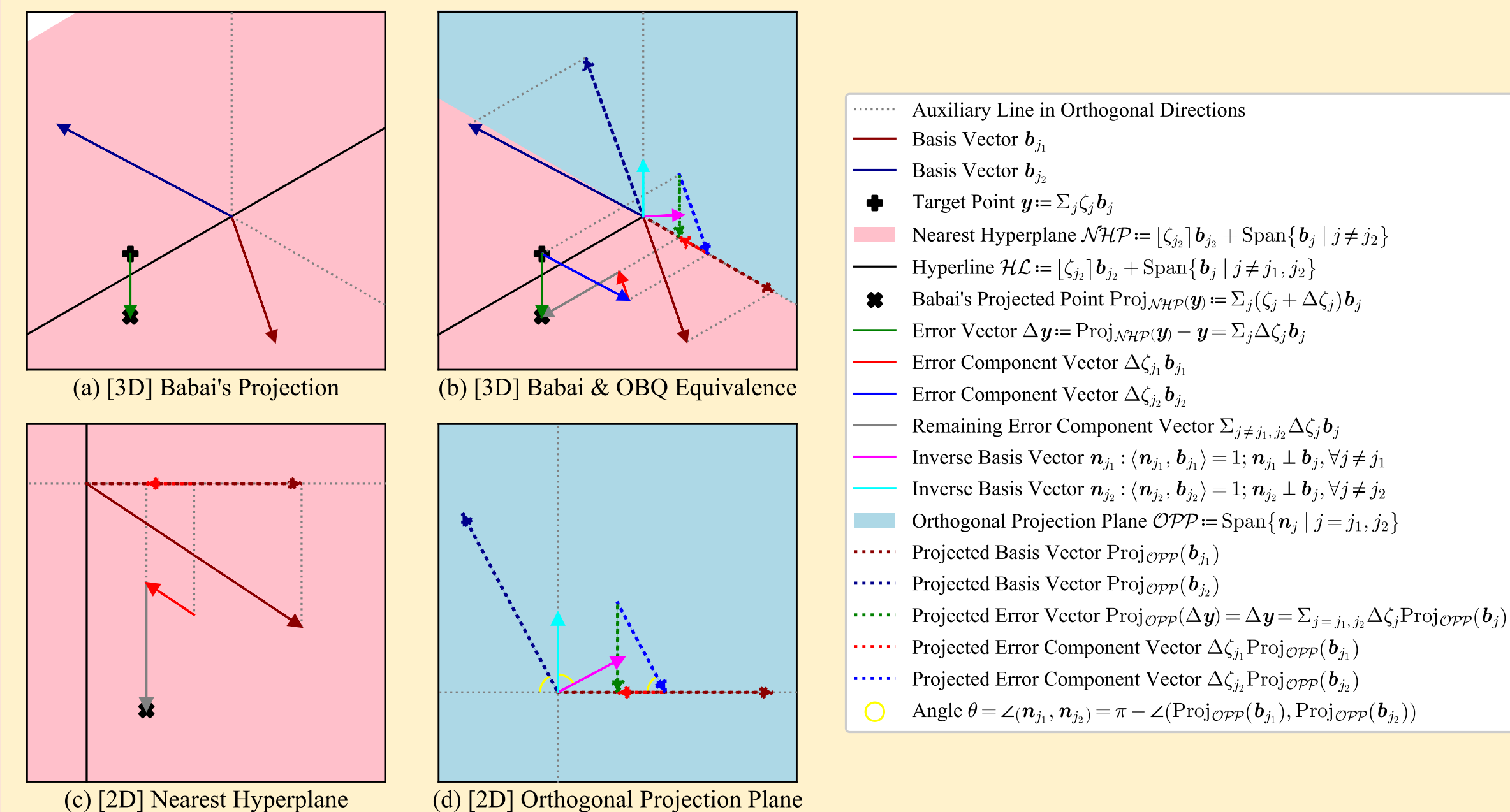


### GPTQ Error Bound & Quantization Order

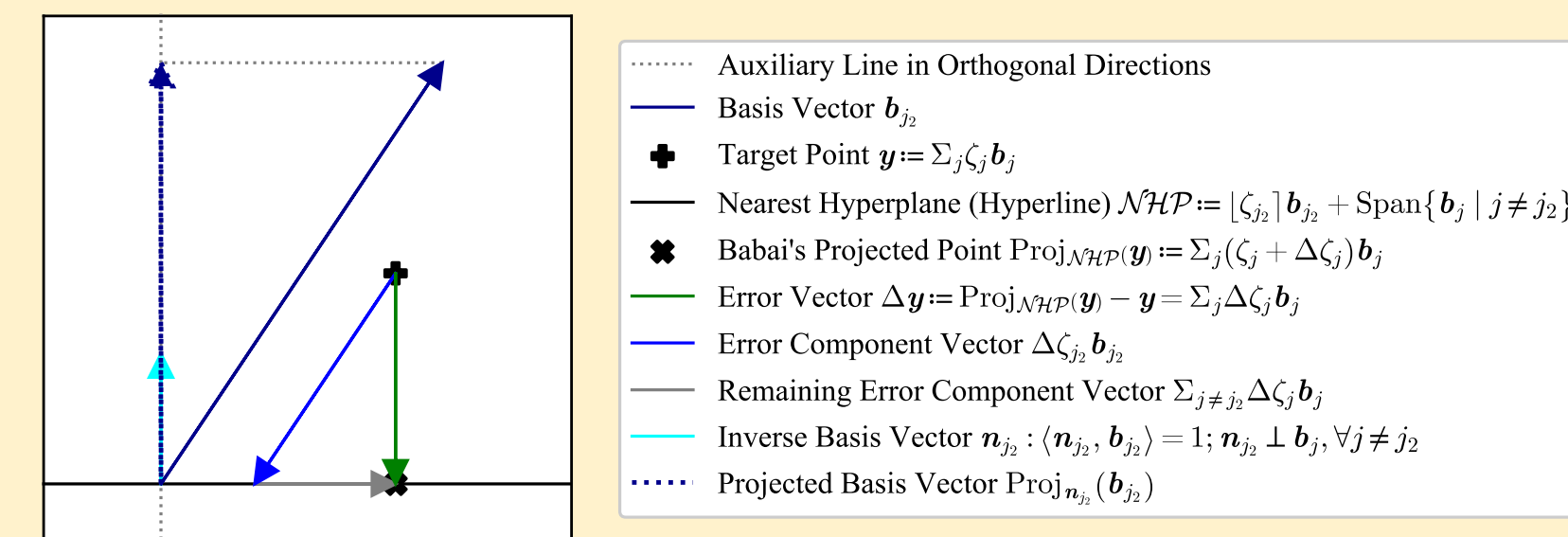


### Proof of Equivalence

#### GPTQ/OBQ Error Propagation $\Leftrightarrow$ Nearest Plane Projection



#### OBQ Quantization Order $\Leftrightarrow$ "Nearest" Nearest Plane



### Applications: Clip-Free GPTQ & Inference Kernel

