

# Comparing Vision Token Interpretability Methods (Layer 0)

**Static NN  
(Input Embed.)**

"Michele"

"blinds"

"rien"

"White"

"handwriting"

**LogitLens  
(LM Head)**

"aber"

"wat"

"Hag"

"marvin"

"wav"

**LN-Lens (Ours)  
(Contextual Embed.)**

"ouver" → "steeple white louvered windows."

L8

"ouver" → "siding of building with louvered ty..."

L8

"ouver" → "a dark green louver"

L8

"ouver" → "kitchen window with black curtain a..."

L8

"-window" → "gabled roof on multi-windowed brow"

L8

**Input Image**

