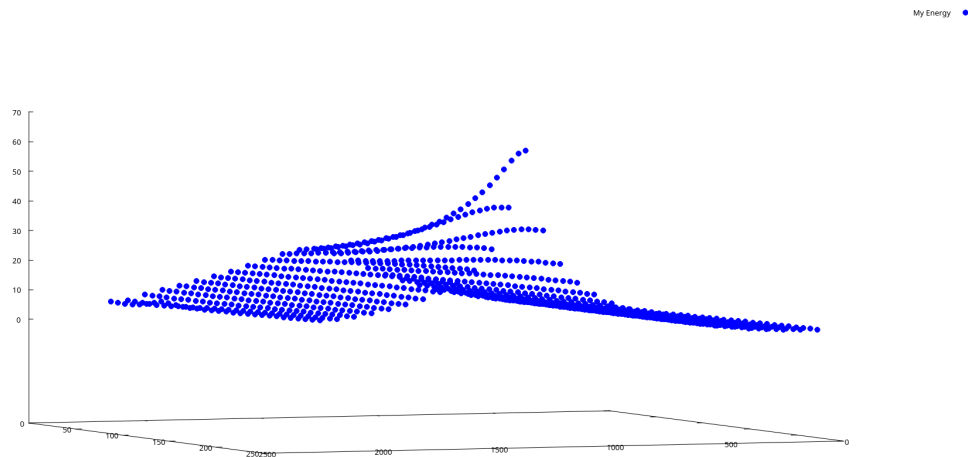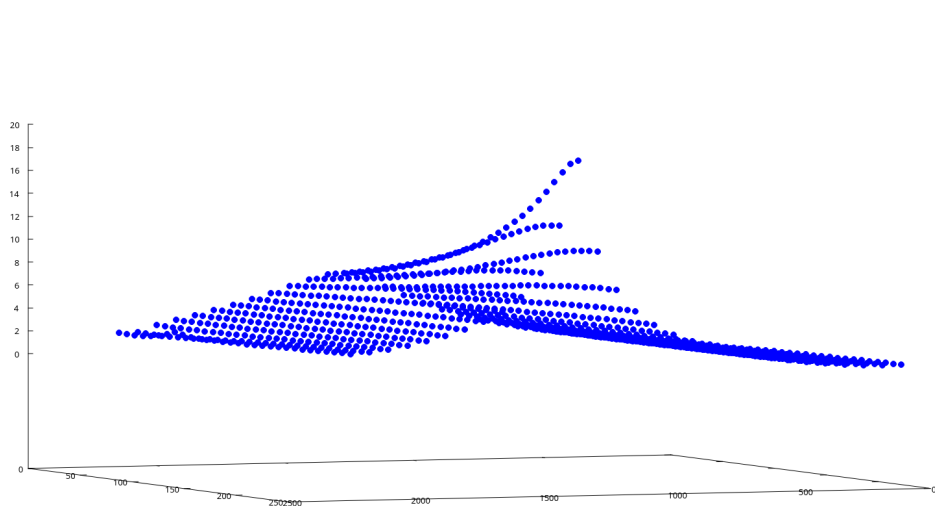I chose to approximate the relationship between kernel configuration parameters (height×width dimension and input channels) and performance metrics (power, energy, and latency) for the add_relu operation across several model configurations.
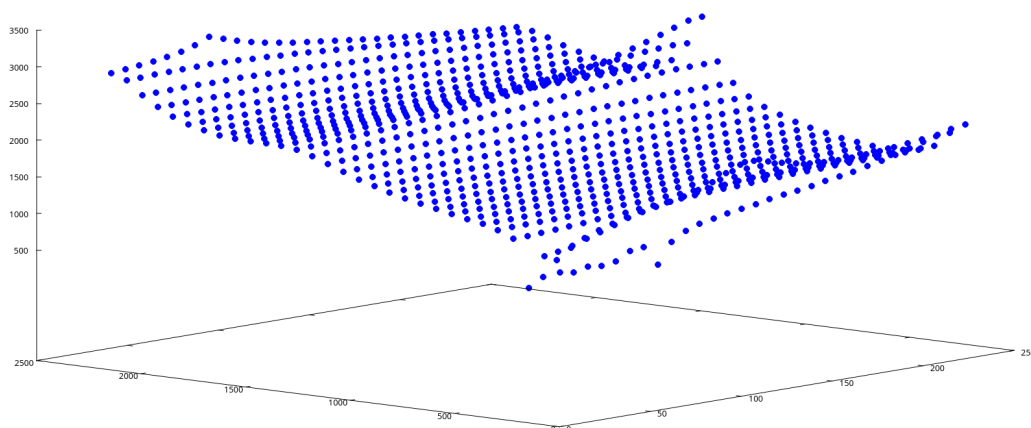
Inputs: Height×Width (HW), Input Channels (CIN)

Outputs: Power (mW), Energy (µJ), Latency (ms)



Energy



Latency

Power

Based on initial data visualization, the 3D plots demonstrated a linear relationship with positive correlation between HW, CIN, and all three output metrics. There seemed to be no need to use a complicated model like MLP or CNN as there is no spatial relationship between the data. Linear regression is well-suited for this application because it is computationally efficient, requires little memory, and is easy to debug. Linear regression has minimal training overhead. The model trains via a closed-form solution (the normal equation: $\beta = (X^T X)^{-1} X^T Y$), requiring only matrix operations rather than iterative optimization. Training time is effectively instantaneous even on CPU. Once trained, the model stores only a small coefficient matrix (3×3 in this case) - approximately 9x8=72 bytes. Inference requires a single matrix multiplication, making it ideal for resource-constrained or real-time scenarios. There is no back-propagation involved. Linear regression finds the best-fit linear relationship between inputs and outputs by minimizing the sum of squared errors. For our multi-input, multi-output case, it learns three separate equations:

Power $= \beta_0^{(1)} + \beta_1^{(1)} \cdot HW + \beta_2^{(1)} \cdot CIN$

Energy $= \beta_0^{(2)} + \beta_1^{(2)} \cdot HW + \beta_2^{(2)} \cdot CIN$

Latency $= \beta_0^{(3)} + \beta_1^{(3)} \cdot HW + \beta_2^{(3)} \cdot CIN$

The algorithm constructs a design matrix X with our input features plus an intercept term, and solves for coefficients that best predict the output matrix Y. The normal equation guarantees finding the globally optimal solution in one step.

Current MSE: 94,563.9

This high error is primarily due to the mismatch in scale between power, latency, and energy. As you can see from the following output where the rightmost three columns are in order: power, latency, energy:



power is much larger than the other values.

For the following lab I will work on the following:

Feature normalization: Standardize inputs and outputs to comparable scales (z-score normalization or min-max scaling)

Proper train/test split: Reserve 20-30% of data for unbiased evaluation

Expanded dataset: Include more samples from add_relu and other operations

The linear model provides a strong baseline with negligible overhead, making it a suitable starting point for kernel performance prediction.