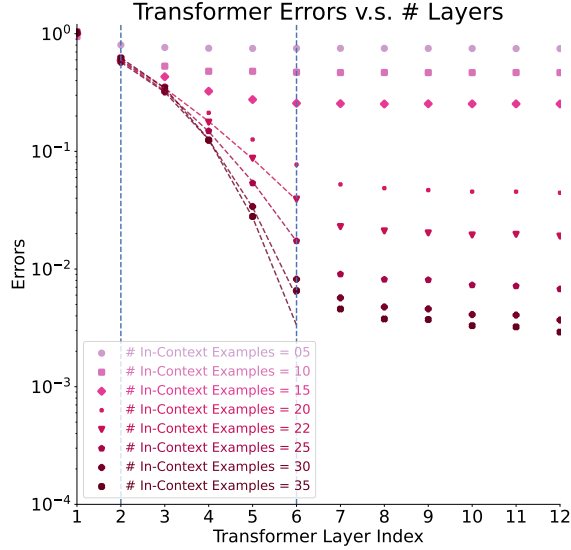# 1 Noisy Linear Regression

We repeat the same experiments on noisy linear regression tasks with $y = w^\top x + \varepsilon$ where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ with noise level $\sigma = 0.1$. As shown in Figure 1a, Transformers still show superlinear convergence on noisy linear regression tasks. Since the predictor is $\hat{w} = \left(X^\top X + \lambda I\right)^+ X^\top y$ for some $\lambda$, the iterative newton's method is applied to $S = X^\top X + \lambda I$. Iterative Newton's method still keeps the same superlinear convergence rates. As shown in Figure 1b, Transformers and Iterative Newton's rates match linearly, as in the noiseless linear regression tasks.



(a) Transformers have superlinear convergence rate

(b) Transformers match Iterative Newton's rate

Figure 1: Empirical Results on Noisy Linear Regression

# 2 Effects of Hidden Dimension

We first restate the main theorem here, with text highlighted in red related to the hidden dimension.

**Theorem 5.1.** *There exist Transformer weights such that on any set of in-context examples $\{x_i, y_i\}_{i=1}^n$ and test point $x_{\text{test}}$, the Transformer predicts on $x_{\text{test}}$ using $x_{\text{test}}^\top \hat{w}_k^{\text{Newton}}$. Here $\hat{w}_k^{\text{Newton}}$ are the Iterative Newton updates given by $\hat{w}_k^{\text{Newton}} = M_k X^\top y$ where $M_j$ is updated as $M_j = 2M_{j-1} - M_{j-1} S M_{j-1}, 1 \le j \le k, M_0 = \alpha S$, for some $\alpha > 0$ and $S = X^\top X$. The dimensionality of the hidden layers is $\mathcal{O}(d)$, and the number of layers is $k + 8$.*

Now we ablate 12-layer 1-head Transformers with various hidden sizes on $d = 20$ problems. As shown in Figure 2, we observe that Transformers can mimic OLS solution when the hidden size is 32 or 64, but fail with smaller sizes. This resonates with our theoretical results on $\mathcal{O}(d)$ hidden dimension, and in this case, Theorem 5.1 ensures a construction of transformers to implement Iterative Newton's method.
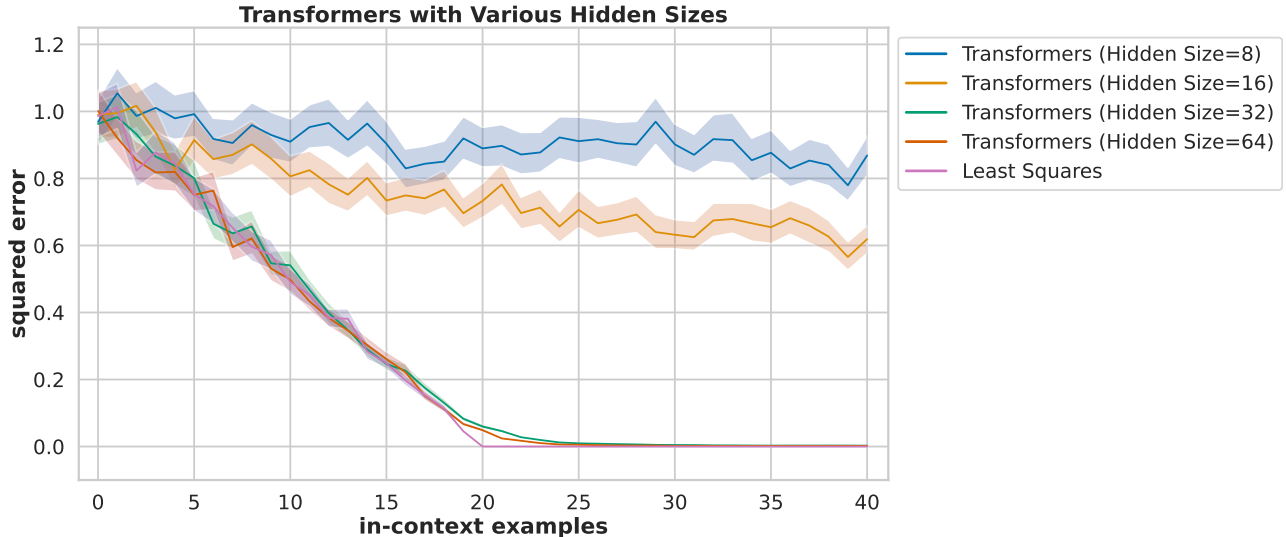


Figure 2: Ablation on Transformer's Hidden Size