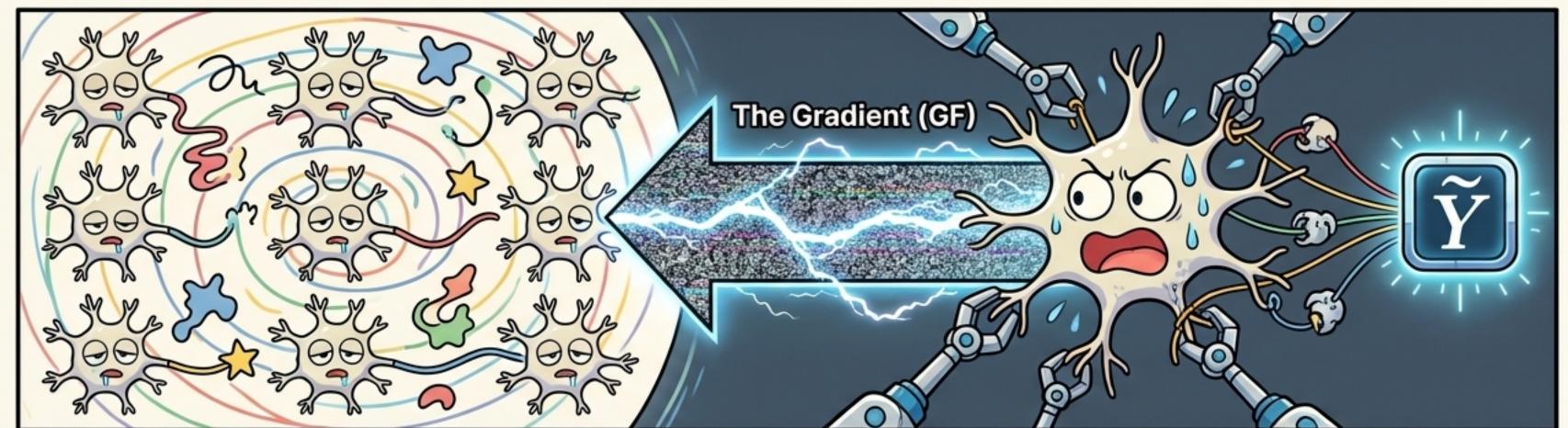
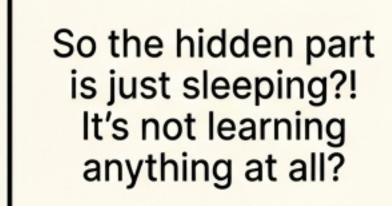


Hidden Layer (W)

Stage I: Lazy Learning

Output Layer (V)





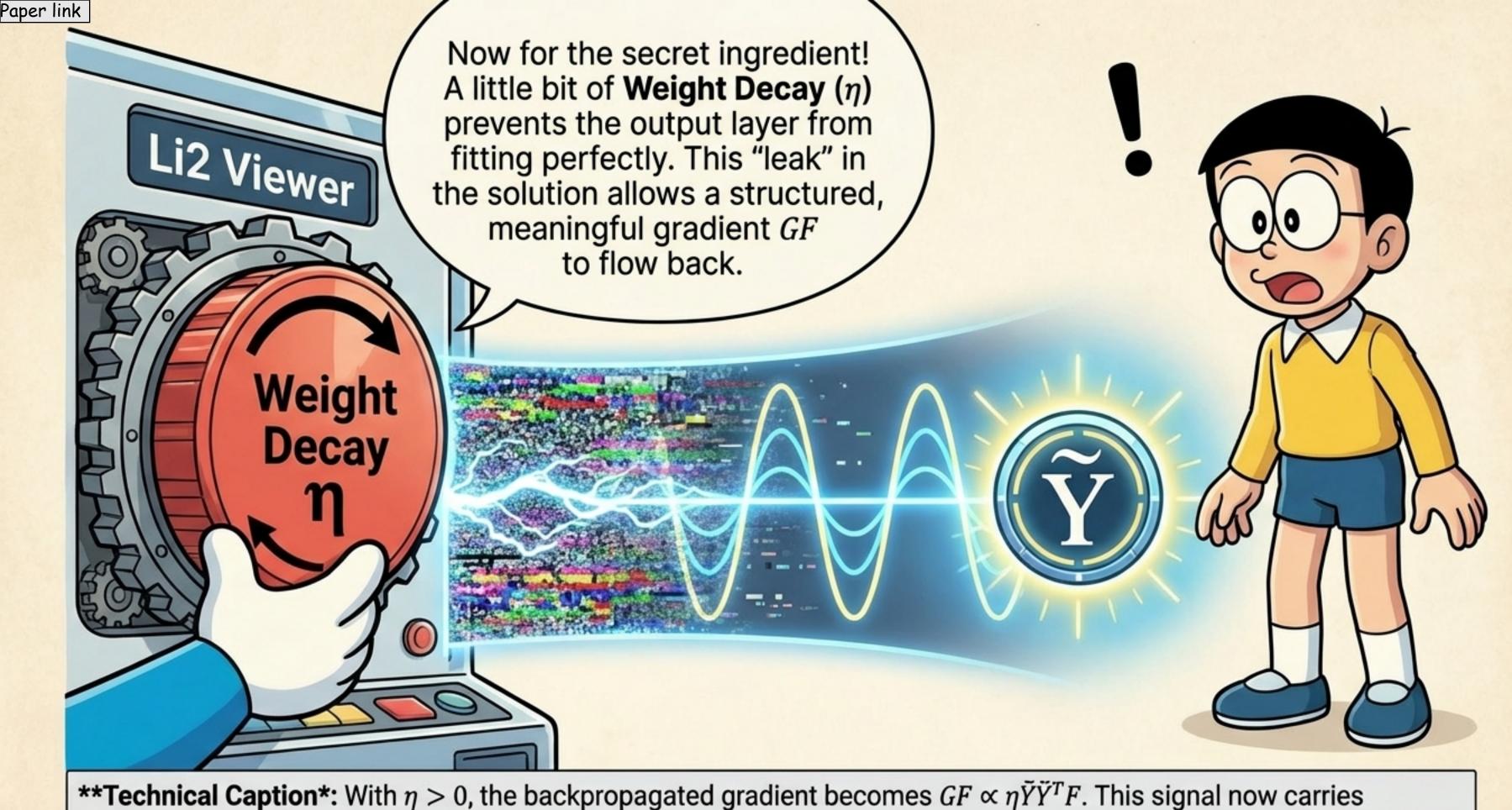


Exactly! In Stage I, only the output layer 'V' learns. It finds a temporary solution—a ridge regression—using the random features 'F' from the hidden layer.

This is the "memorization circuit".

NotebookLM

Technical Caption: The backpropagated gradient $GF = P^{\perp 1}(Y - FV)V^T$ is initially noise. Without weight decay $(\eta = 0)$, GF' eventually becomes zero, and feature learning stops entirely.



**Technical Caption*: With $\eta > 0$, the backpropagated gradient becomes $GF \propto \eta \tilde{Y} \tilde{Y}^T F$. This signal now carries information about the target label, triggering the start of feature learning in the hidden layer. (Source: Lemma 1, Eqn. 6) NotebookLM

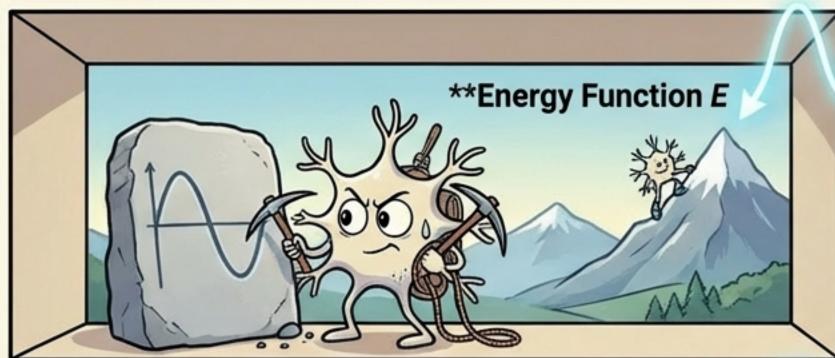
Paper link

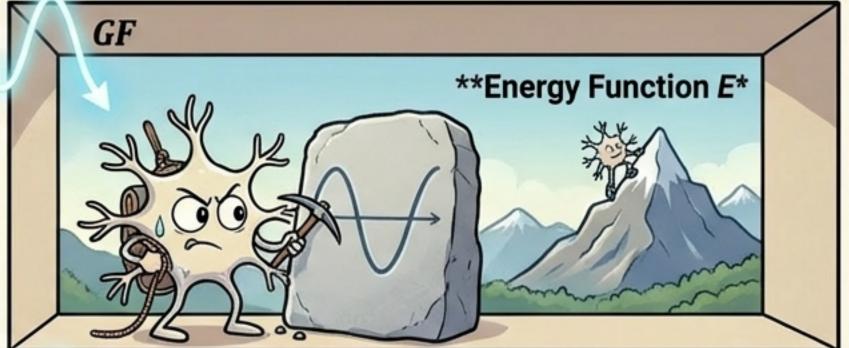


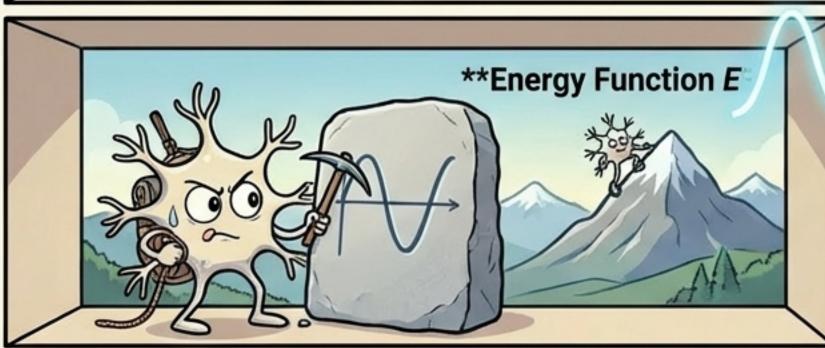
Wow! They all woke up!
And they're learning
structured patterns, but all
by themselves?

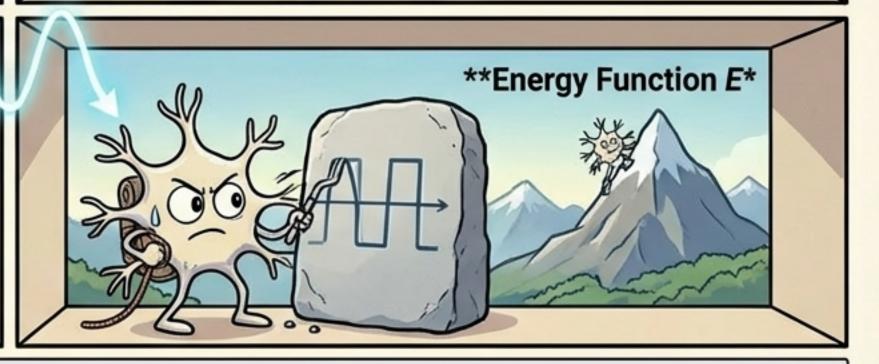
Correct. The gradient *GF* now acts on each neuron independently. Each one climbs its own "Energy Function" *E*". The peaks—the local maxima—are the very features they learn!



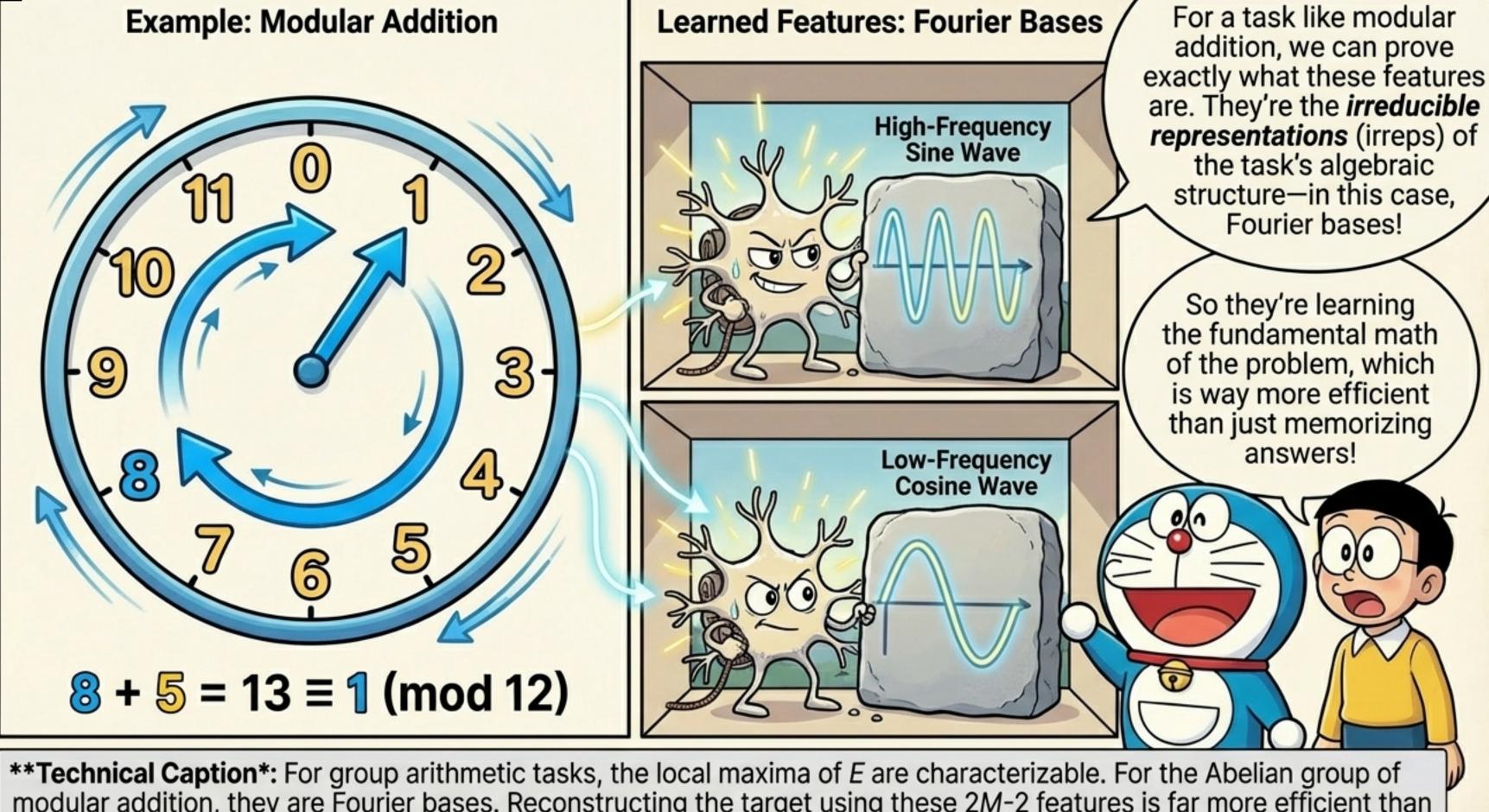






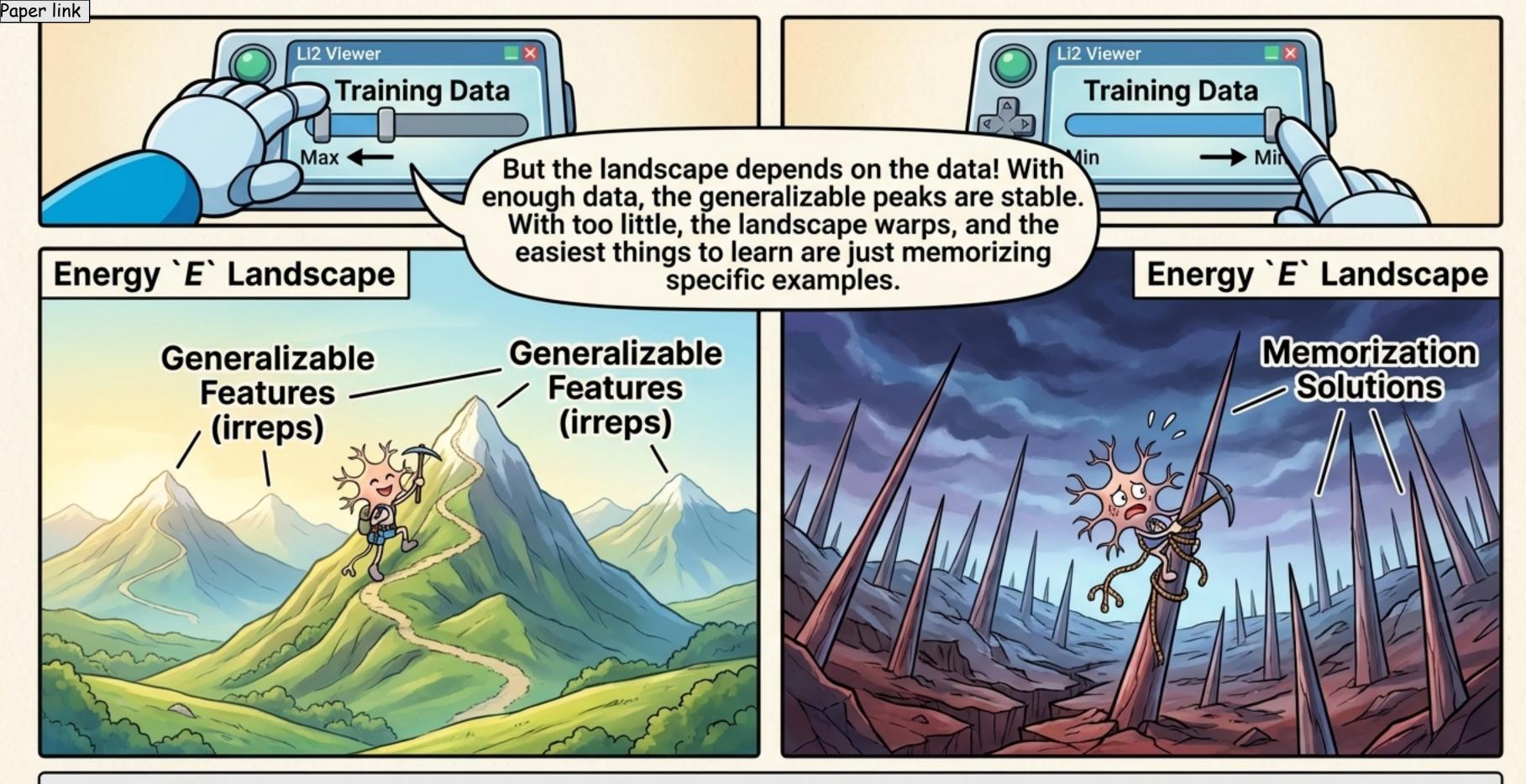


**Technical Caption*: The dynamics of each hidden node w_j is exactly the gradient ascent of an energy function: $E(w_j) = \frac{1}{2} ||\tilde{Y}^T \sigma(Xw_j)||^2$. The emerging features are the local maxima of E. (Source: Theorem 1)



Paper link

**Technical Caption*: For group arithmetic tasks, the local maxima of E are characterizable. For the Abelian group of modular addition, they are Fourier bases. Reconstructing the target using these 2M-2 features is far more efficient than a memorization solution requiring M^2 nodes. (Source: Corollary 2, Theorem 3) NotebookLM



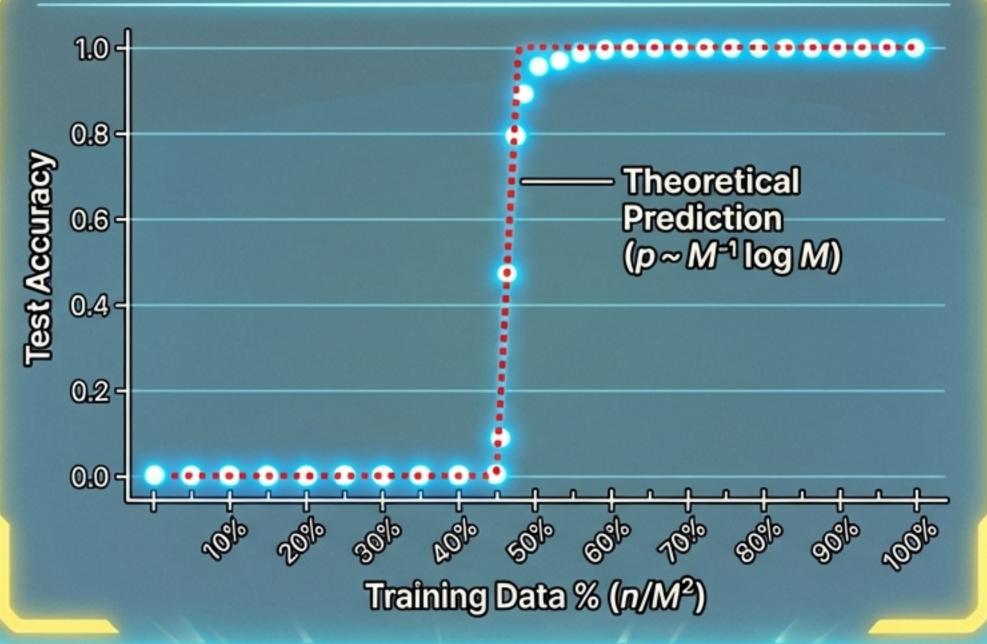
Technical Caption: Data governs the landscape of `E`. Sufficient training data maintains the shape of generalizable local maxima. Insufficient data leads to the emergence of non-generalizable, memorization-based local maxima. (Source: Theorems 4 & 5)

⋒ NotebookLM

This allows us to derive a provable scaling law! We can predict the precise amount of data needed for the network to 'grok'. The theory matches the experiment perfectly.



Theory Meets Experiment: A Provable Scaling Law

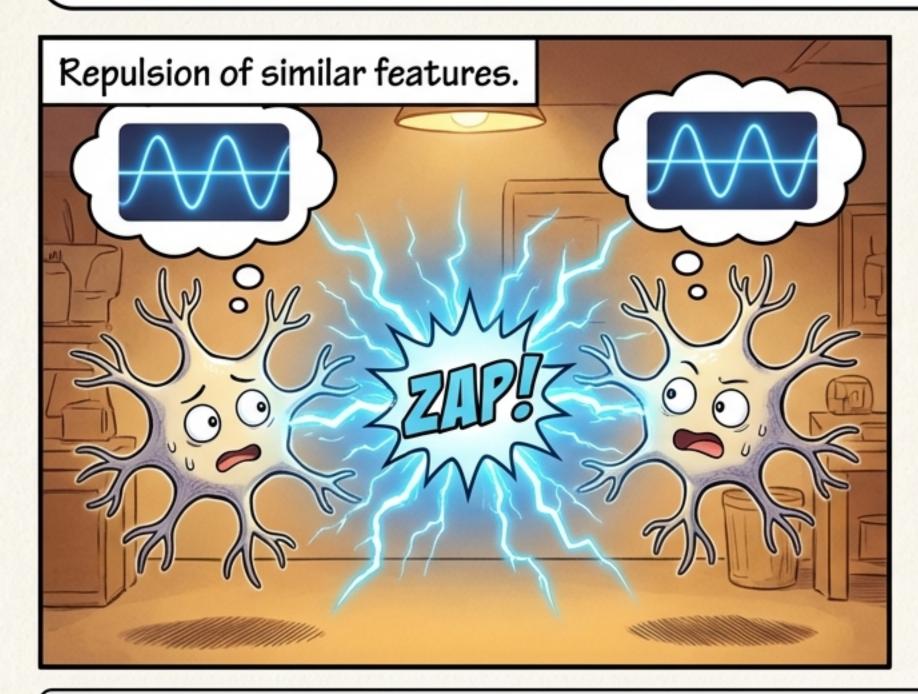


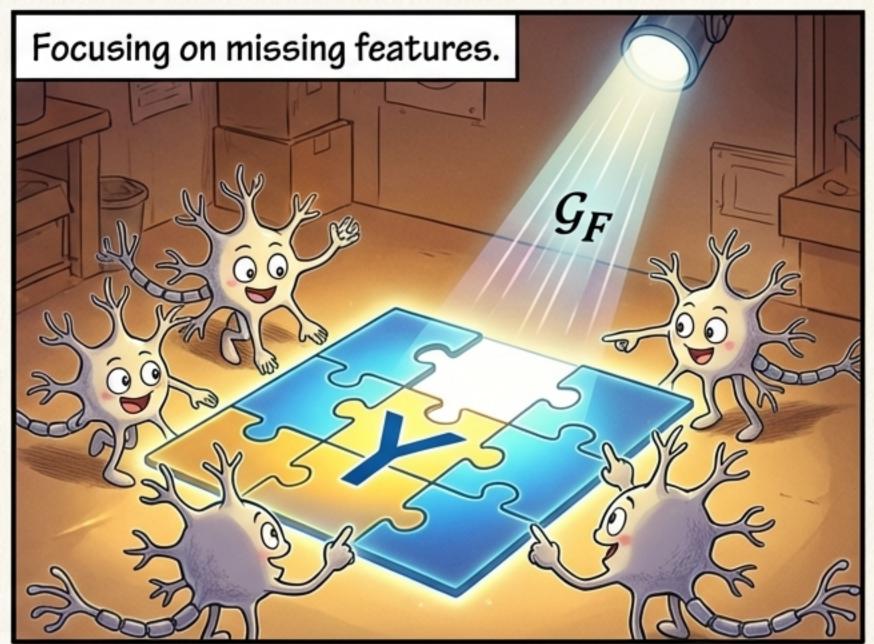
So it's not a mystery at all!
There's a mathematical rule for when the big jump in understanding happens!



In Stage III, the neurons are no longer independent. They start to interact. Neurons with similar features repel each other to promote diversity. The gradient \mathcal{G}_F also changes to guide them, pointing out the parts of the problem they haven't solved yet.

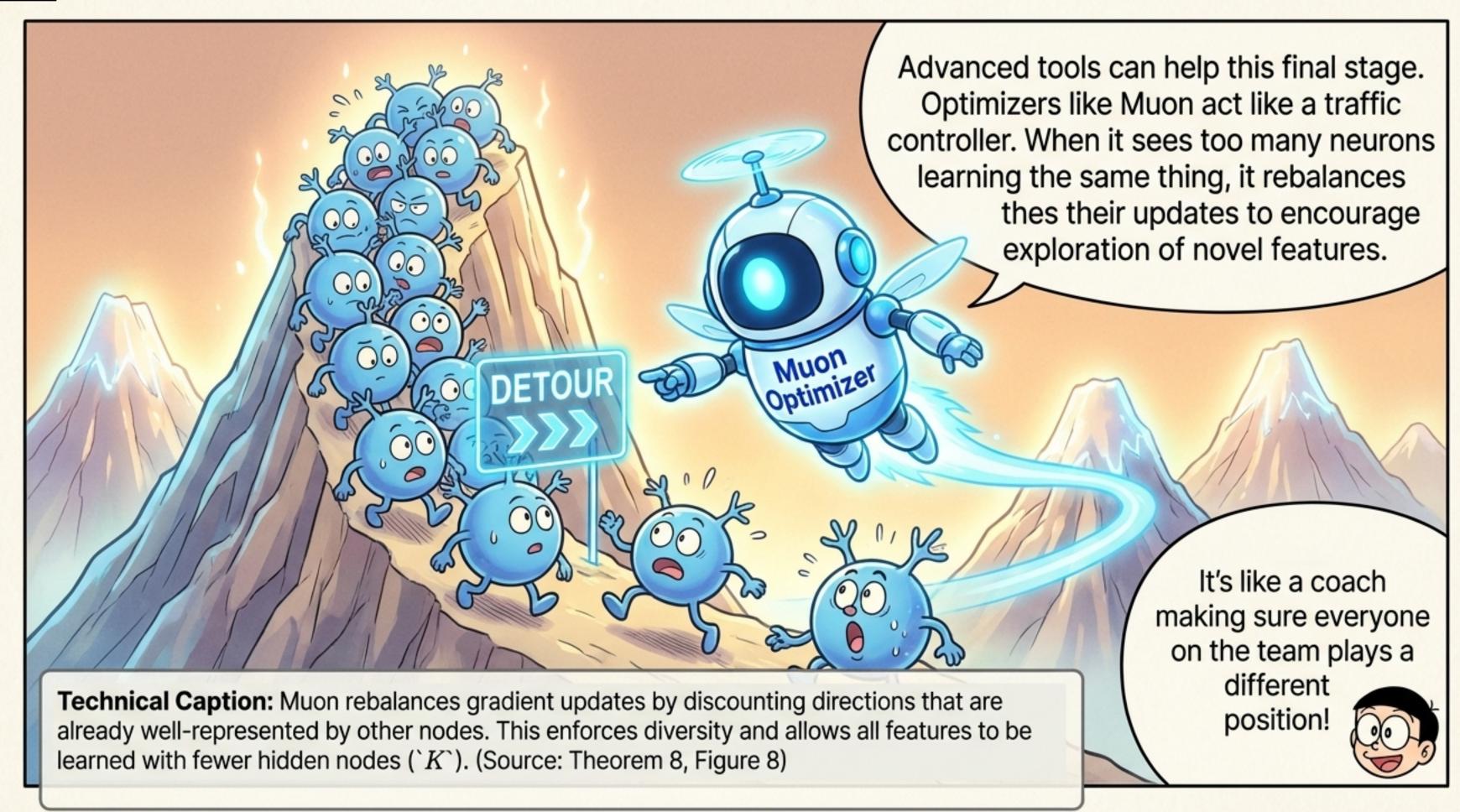


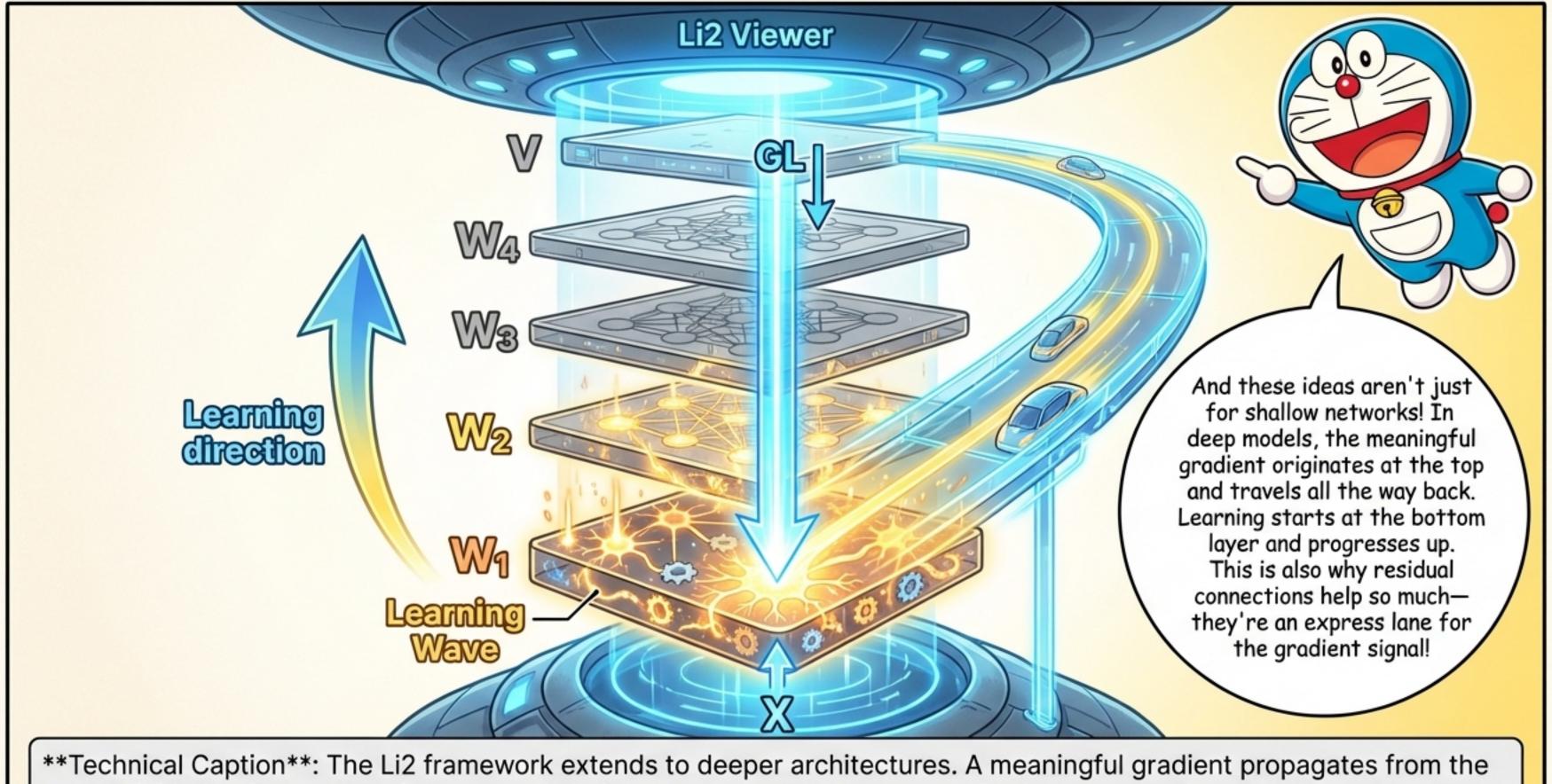




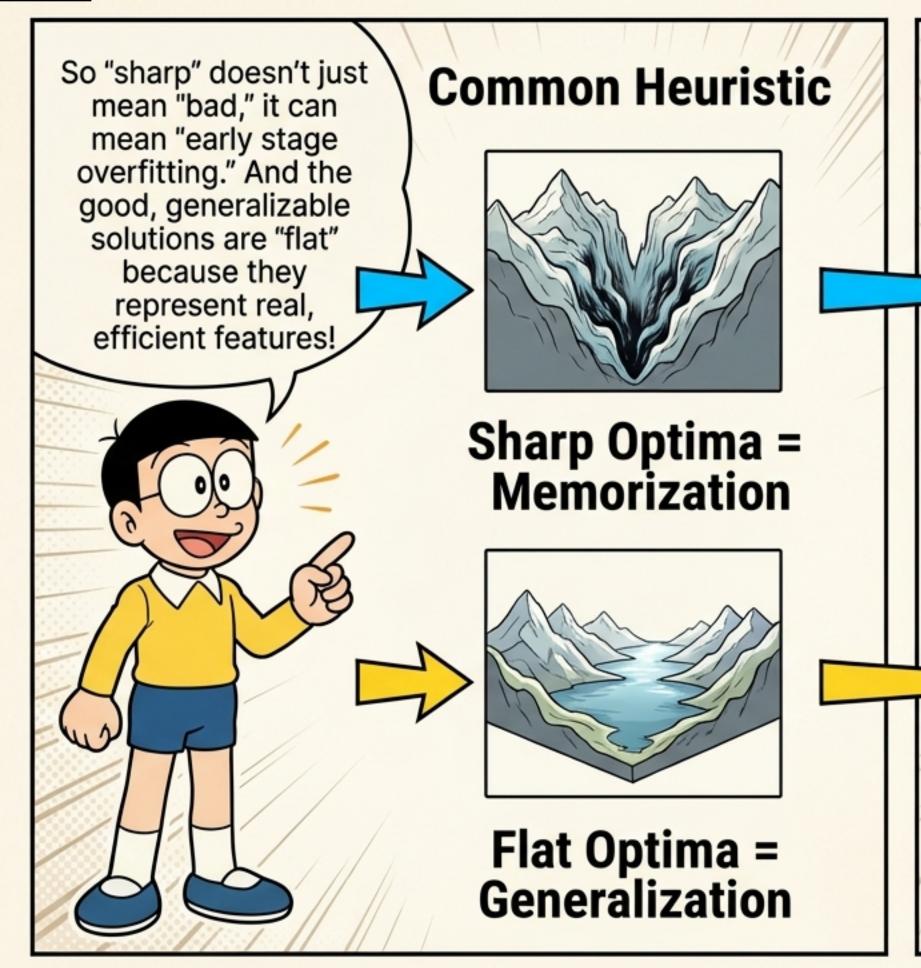
Technical Caption:

- Repulsion (Thm. 6): Off-diagonal terms in the gradient dynamics push nodes with similar activations apart.
- Top-down Modulation (Thm. 7): Once a subset of irreps S is learned, GF changes to create a modified energy function E_S with local maxima only on the missing irreps.





Technical Caption: The Li2 framework extends to deeper architectures. A meaningful gradient propagates from the output, initiating feature learning sequentially from the lowest layer upwards. Residual connections provide a cleaner, stronger signal to these lower layers. (Source: Sec 6)



Li2 Framework View



Stage I: Overfitting on random features



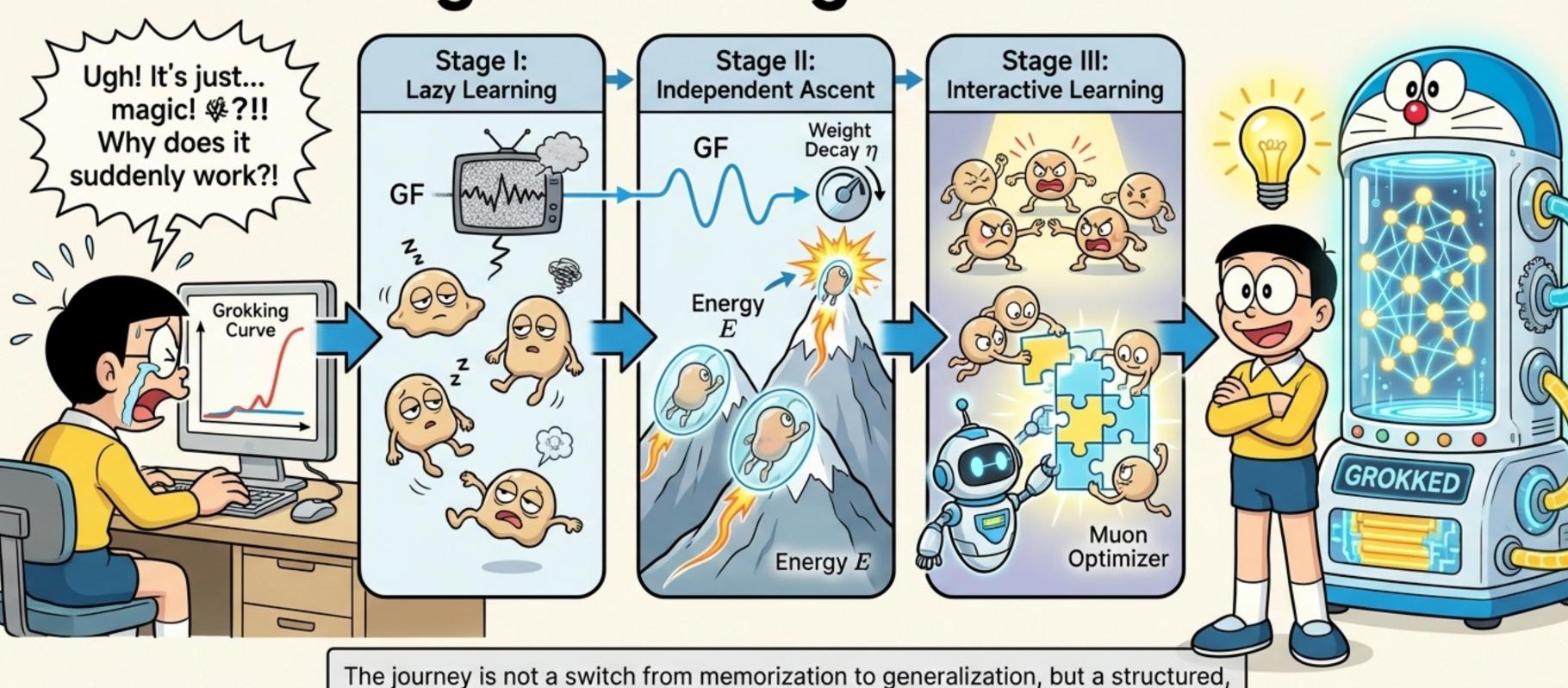
Stage II: Learned features are provably flat local maxima of `E`

memortion in grokking corresponds to the sharp optimum of the Stage I ridge solution.

True generalization comes from finding the provably flat local maxima of the energy function E'. (Source: Sec 7, Corollary 1)

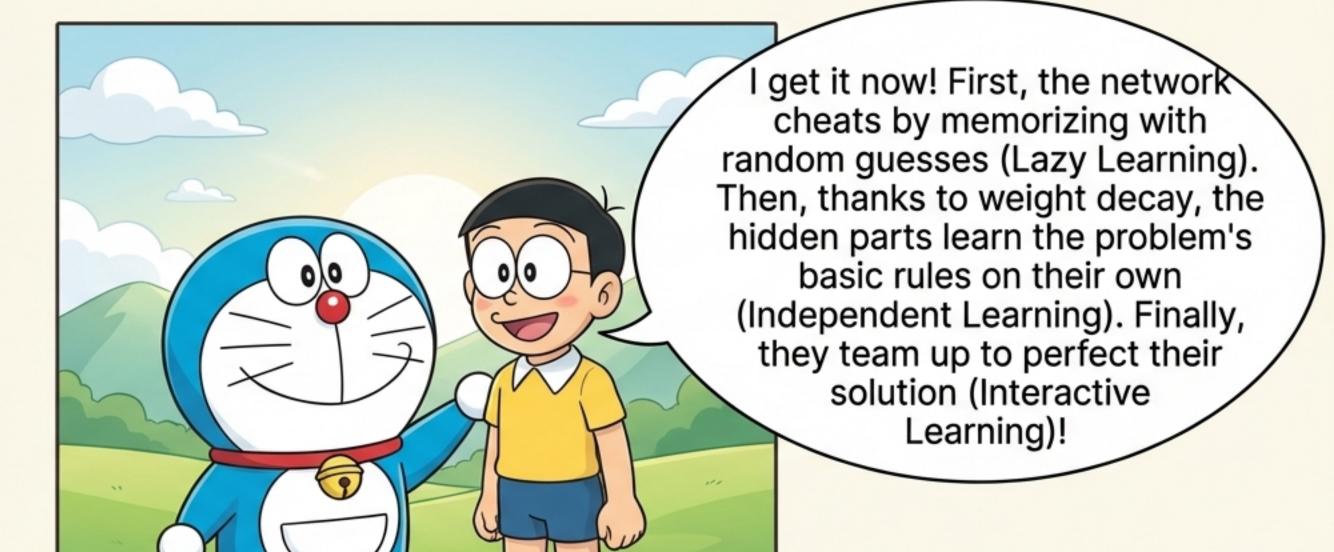


Grokking: From Magic to Mechanism.



three-stage process of feature emergence driven by gradient dynamics.

€ NotebookLM



Doraemon's Concluding Thought

The Li2 framework demystifies grokking by modeling the precise dynamics of how features emerge. It reveals that grokking is **the observable effect of a network transitioning from overfitting on random features** to **learning the true, efficient structure of the data**, governed by a data-dependent energy landscape. By understanding the mechanism, we derive provable scaling laws that connect theory directly to practice.