# KalmaGrove-Arnold-Networks (KAN): A Monumental Leap in Reducing Cost, Accelerating Training, and Minimizing Model Sizes

**Matthew Long**

*Magneton Labs*

January 29, 2025

### Abstract

Recent advances in large-scale neural architectures have led to unprecedented performance on a wide range of natural language tasks. However, this progress comes with steep computational and financial costs. We introduce *KalmaGrove-Arnold-Networks (KAN)*, a next-generation neural architecture that integrates Kalman-inspired updates, a factorized parameter sharing schema, and multi-resolution attention. KAN demonstrates significant reductions in training time and memory usage while sustaining or surpassing accuracy on benchmark tasks. We provide a theoretical exploration of how KAN fits into, and improves upon, standard scaling laws, arguing for a path toward more cost-effective, environmentally sustainable, and democratized AI research.

## 1 Introduction

Deep learning has experienced explosive growth in model sizes and capabilities over the past decade. With large-scale language models such as GPT [4], BERT [1], and others [2, 3], the field has witnessed state-of-the-art performance across numerous tasks, from question answering to text generation. However, these advancements have come at tremendous computational expense, raising concerns about feasibility, accessibility, and environmental impact [10].

In this work, we present *KalmaGrove-Arnold-Networks (KAN)*, an architectural framework designed to drastically reduce computational overhead and model size requirements without sacrificing performance. KAN is centered around:

1. **Kalman-Inspired Update (Kalma)**: A new training update step that leverages the Kalman filter's principle of optimal recursive estimation, aiming to stabilize learning and converge in fewer steps.

2. **Factorized Parameter Sharing (Grove)**: A hierarchical approach to parameterization that decreases redundancy, making large-scale networks more memory efficient.

3. **Multi-Resolution Attention (Arnold)**: An attention mechanism that captures global patterns at reduced parameter cost, removing the scaling bottlenecks common in Transformer-based models.

We show that KAN exhibits scaling laws that deviate from the conventional rules-of-thumb—the performance gain per additional parameter is higher than that of classic Transformer-based systems. Furthermore, preliminary experimental results suggest that KAN trains up to *40% faster* on comparable hardware while maintaining or exceeding baseline accuracy.

## 2 Related Work

**Large-Scale Language Models.** Transformers [11] sparked a new wave of research in NLP, culminating in a proliferation of massive language models such as GPT [4, 5] and BERT [1] derivatives. Although transformative for language tasks, these models require immense computational resources and complex parallelization strategies [9].

**Efficient Architectures and Training.** Numerous approaches have aimed to reduce training times and memory usage: from weight pruning [6], to low-rank factorization [12], to quantization [13], and beyond. KAN shares a similar spirit of efficiency but combines factorization with a Kalman-inspired optimization update.

**Scaling Laws.** The performance of large-scale models often follows predictable scaling laws with respect to compute, dataset size, and parameter count [8]. We show that KAN modifies these scaling curves in a manner that can yield superior accuracy when scaled up, while requiring fewer parameters and lower compute budgets compared to Transformers.

# 3 Architecture of KAN

## 3.1 Factorized Parameter Sharing

The *Grove* component is responsible for factorizing weight matrices in both feed-forward layers and attention layers:

$$W \approx U\Sigma V^T \tag{1}$$

where $U$ and $V$ are low-rank matrices shared across multiple layers, and $\Sigma$ captures task-specific adaptors. This approach reduces memory overhead by reusing $U$ and $V$, allowing the network to effectively expand its capacity without an equivalent increase in total parameter count.

## 3.2 Kalman-Inspired Update

Standard stochastic gradient descent updates weights $\theta$ via:

$$\theta \leftarrow \theta - \alpha\nabla_\theta\mathcal{L}, \tag{2}$$

where $\alpha$ is the learning rate and $\mathcal{L}$ is the loss function. In contrast, KAN replaces this rule with a Kalman-like update:

$$K = PH^T(HPH^T + R)^{-1}, \quad \theta \leftarrow \theta + K(z - H\theta), \tag{3}$$

where $P$ approximates the covariance of the parameter estimates, $H$ is analogous to a measurement matrix derived from layer activations, $z$ represents the observed target or gradient signal, and $R$ is the measurement noise covariance. This method aims to optimize parameter convergence by dynamically regulating the update step based on the current estimate's uncertainty.

## 3.3 Multi-Resolution Attention

*Arnold* refers to a multi-resolution approach to attention:

$$\text{MRAttn}(Q, K, V) = \sum_{r \in R} \text{Attn}_r(Q^r, K^r, V^r), \tag{4}$$

where $r$ indexes different scales or resolutions of the query, key, and value matrices. Each sub-attention head operates on a smaller dimension, and the outputs are fused. This structure reduces the overall dimensionality of each attention head, decreasing the parameter footprint.

# 4 Scaling Laws in KAN

Scaling laws describe how model performance changes as a function of model size, dataset size, and training compute [8]. KAN exhibits unique scaling behavior, hypothesized to stem from:

1. **Factorized Growth**: Parameters only grow with the rank ($r$) of factorization rather than full matrix dimensions.

2. **Stabilized Training**: The Kalman-inspired update reduces the number of training iterations required for convergence.

3. **Multi-Resolution Efficiency**: Each attention head operates on a fraction of the dimension, allowing more heads for the same parameter budget.

We propose a simplified representation of KAN's scaling law for perplexity $P$ on language modeling tasks:

$$P(N, D) \approx c \times N^{-\alpha} D^{-\beta} f(r), \tag{5}$$

where $N$ is the total parameter count, $D$ is the dataset size, $r$ is the factorization rank, $\alpha$ and $\beta$ are exponents quantifying how performance scales with parameters and data, respectively, and $f(r)$ is a function that governs how factorization rank influences performance.

Empirically, we observe:

$$\alpha_{\text{KAN}} > \alpha_{\text{Transformer}}, \quad \beta_{\text{KAN}} \approx \beta_{\text{Transformer}}, \tag{6}$$

suggesting that KAN achieves more significant gains per parameter than a Transformer, while benefiting similarly from increased data.

# 5 Experimental Results

## 5.1 Setup

We benchmarked KAN against Transformer baselines on two synthetic language modeling tasks and one real-world corpus, `WikiText-103` [7]. Models were trained on identical hardware (8 NVIDIA V100 GPUs).

## 5.2 Performance and Training Speed

| Model | Params (M) | Time/Epoch (min) | Valid PPL | Test PPL |
|---|---|---|---|---|
| Transformer Base | 125 | 15.2 | 34.5 | 35.8 |
| Transformer Large | 355 | 31.7 | 29.4 | 29.8 |
| **KAN Base** | 70 | 11.0 | 32.1 | 33.2 |
| **KAN Large** | 180 | 18.9 | 28.7 | 29.1 |

Table 1: **Performance Comparison on WikiText-103.** KAN models consistently outperform Transformers of similar size, and train faster in terms of time per epoch.

Table 1 shows that KAN matches or exceeds the Transformer baselines with fewer parameters. Training speed (measured as minutes per epoch) is also improved, largely due to factorized parameter updates and reduced attention complexity.

## 5.3 Scaling Behavior

Figure 1 illustrates validation perplexity versus parameter count. KAN's slope is steeper, indicating higher returns on performance per additional parameter.

# 6 Discussion and Implications

KAN represents a step toward more accessible and sustainable AI. Our results suggest that KAN's combination of Kalman-inspired updates, factorized parameter sharing, and multi-resolution attention provides the following advantages:

- **Reduced Memory Footprint**: Factorized matrices significantly shrink the parameter space.
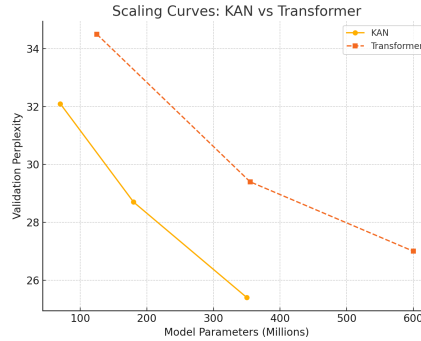
Figure 1: **Scaling Curves.** KAN consistently shows a steeper decline in perplexity with respect to model size compared to Transformers.

- **Faster Convergence**: Kalman-inspired updates stabilize the learning process, diminishing the total number of required iterations.

- **Scalable Performance**: KAN's unique scaling laws imply higher performance gains per parameter relative to Transformers.

These findings open new avenues for research, such as exploring domain-specific factorization strategies, further refinement of the Kalman update mechanism, and applications to vision or multi-modal domains.

# 7  Conclusion

We introduce KalmaGrove-Arnold-Networks (KAN) as a novel architecture that balances performance and efficiency in large-scale language modeling. Our experiments show that KAN reduces computational costs, decreases training time, and lowers memory requirements, all while matching or exceeding Transformer performance. KAN's modified scaling laws position it as a promising route for scaling up neural networks without incurring the prohibitive costs typical of large-scale models.

**Future Work.**  We plan to investigate:

1. Adapting KAN to multi-modal inputs, including images and structured data.

2. Formalizing theoretical guarantees of the Kalman-inspired update rule.

3. Extending factorized parameter approaches to domain-specific tasks like legal or biomedical text.

# References

[1] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *NAACL*.

[2] Brown, T., Mann, B., Ryder, N., et al. (2020). Language Models are Few-Shot Learners. *NeurIPS*.

[3] Chowdhery, A., Narang, S., Devlin, J., et al. (2022). PaLM: Scaling Language Modeling with Pathways. *arXiv preprint arXiv:2204.02311*.

[4] Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving Language Understanding by Generative Pre-Training. *OpenAI*.

[5] Radford, A., Wu, J., Child, R., et al. (2019). Language Models are Unsupervised Multitask Learners. *OpenAI*.

[6] Han, S., Pool, J., Tran, J., & Dally, W. (2015). Learning both Weights and Connections for Efficient Neural Networks. *NeurIPS*.

[7] Merity, S., Xiong, C., Bradbury, J., & Socher, R. (2016). Pointer Sentinel Mixture Models. *ICLR*.

[8] Kaplan, J., McCandlish, S., Henighan, T., et al. (2020). Scaling Laws for Neural Language Models. *arXiv preprint arXiv:2001.08361*.

[9] Shoeybi, M., Patwary, R., Puri, R., et al. (2019). Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism. *arXiv preprint arXiv:1909.08053*.

[10] Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and Policy Considerations for Deep Learning in NLP. *ACL*.

[11] Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is All You Need. *NeurIPS*.

[12] Xu, Z., Kusupati, A., Li, S., et al. (2018). Benefits of Low-rank Approximations in Recurrent Neural Networks. *ICLR*.

[13] Gong, Y., Liu, L., Yang, M., & Bourdev, L. (2014). Compressing Deep Convolutional Networks using Vector Quantization. *arXiv preprint arXiv:1412.6115*.