



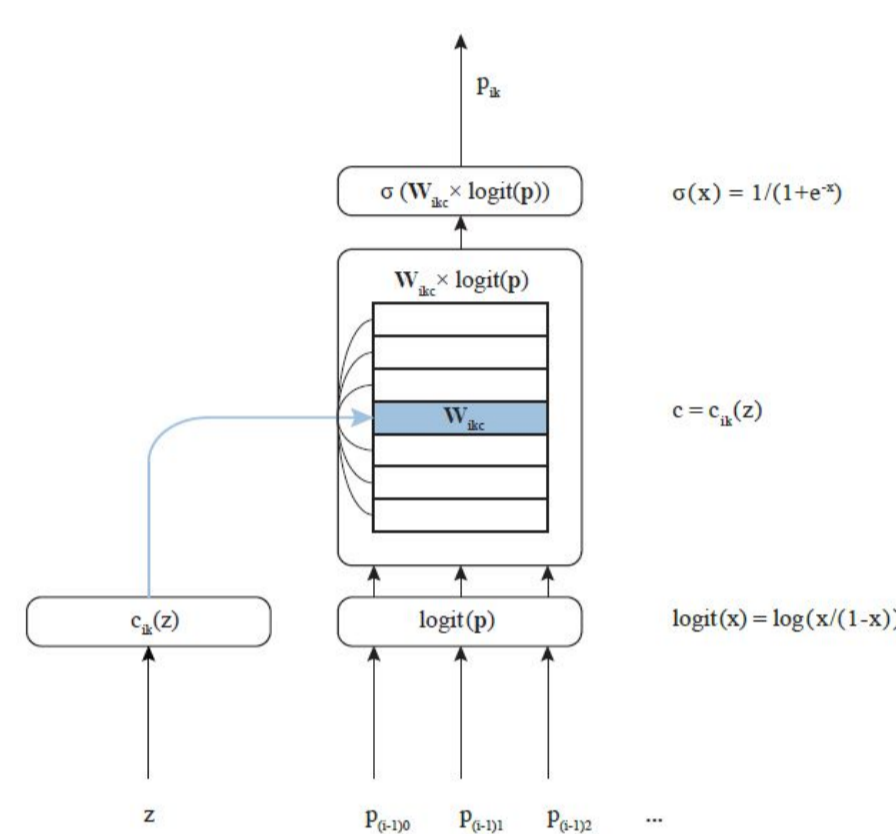
# DeepMind Gated Linear Networks

Joel Veness\*, Tor Lattimore\*, David Budden\*, Avishkar Bhoopchand\*, Christopher Mattern, Agnieszka Grabska-Barwinska, Eren Sezener, Jianan Wang, Peter Toth, Simon Schmitt, Marcus Hutter

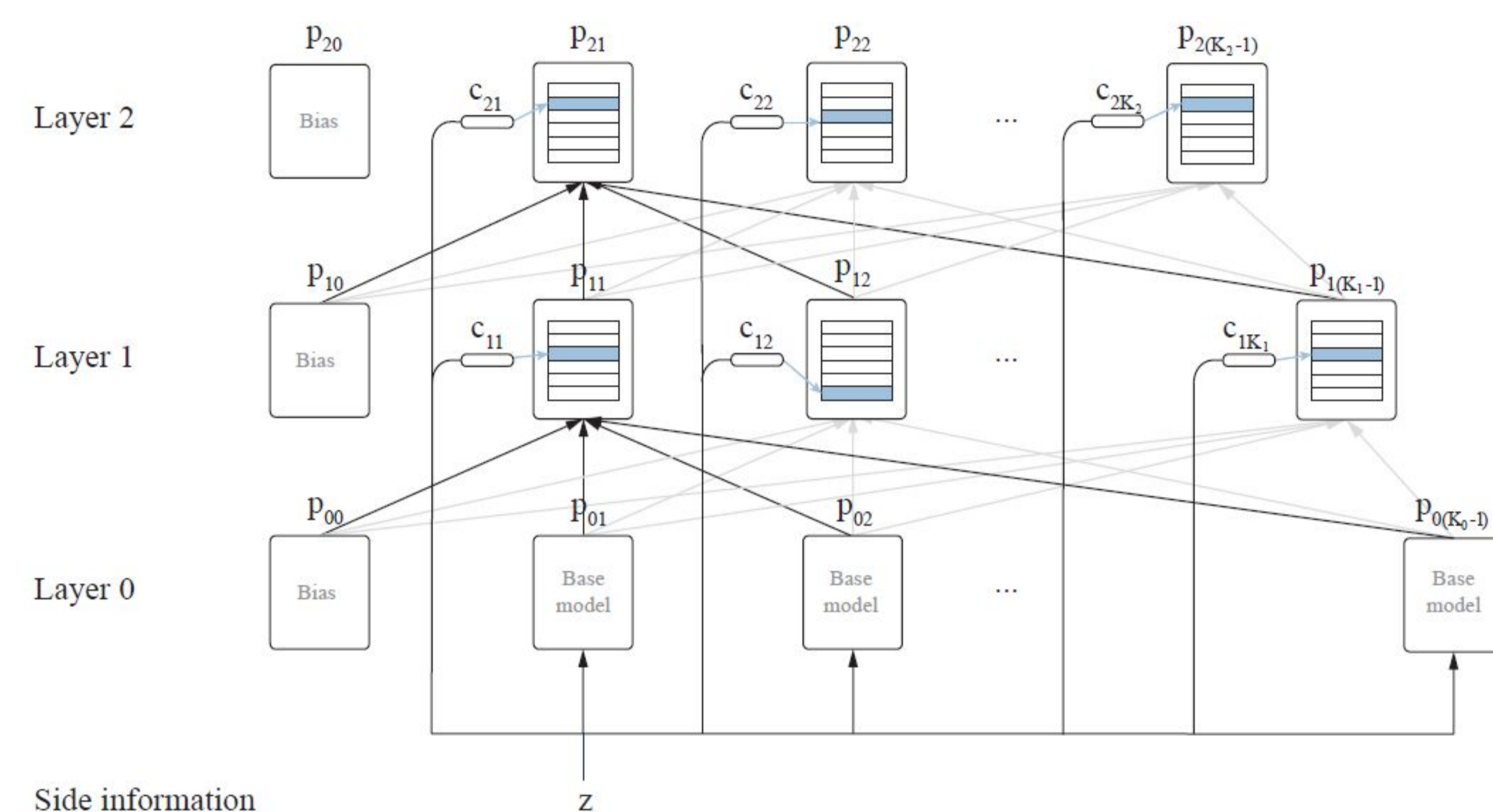
## Introduction

- Gated Linear Networks (GLNs) are a general purpose family of neural networks, with an interesting and distinct take on credit assignment.
- Many possible practical uses, such as regression [1], contextual bandits [3], multiclass classification [4], transfer learning and non-stationary time series modelling [2], density estimation [5] and data compression.

- Associated to each neuron is a fixed, deterministic context function which given an input selects a weight vector to use
- Log loss w.r.t. local weights is convex => can use OGD for optimisation.
- Input/output non-linearities cancel.



## GLN Architecture



## Key Properties

- Fast, online learners
- Robust to forgetting
- Universal
- Parameters adapted via convex optimisation
- Can think of each neuron acting as a kind of distributional "boosting" mechanism
- The loss functions are **convex** for exponential family PDFs.
- Targets and input are broadcast to each neuron, each neuron has own loss function.

- Interpretable
- Distributional
- Compositional



## Gating + Local Learning => Universality

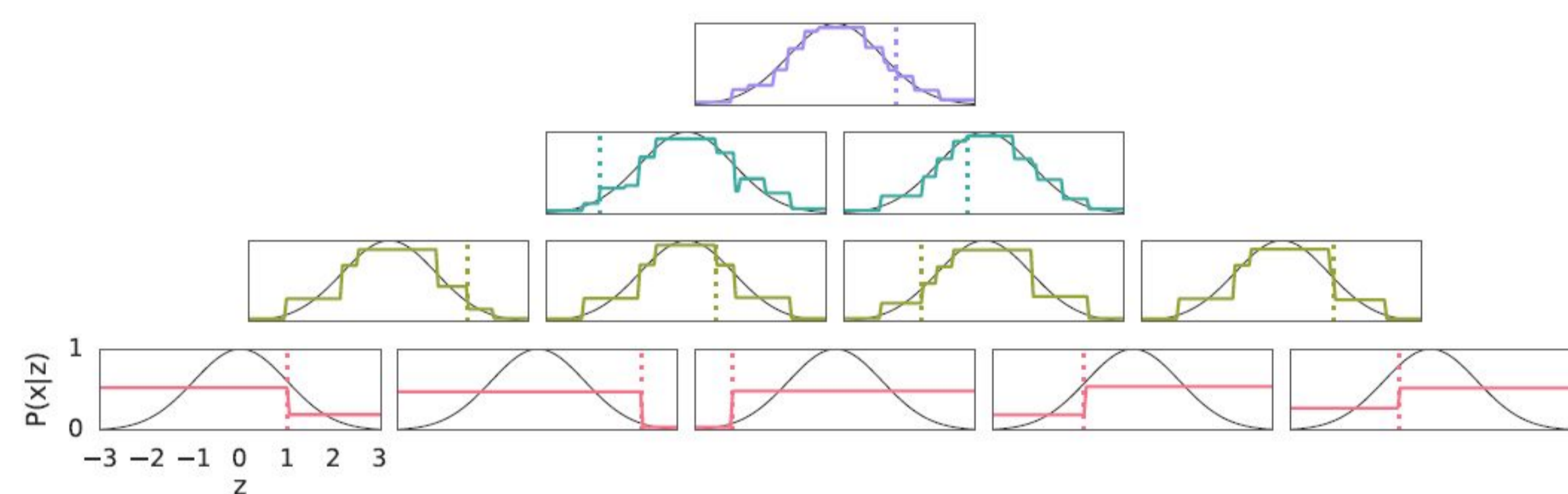
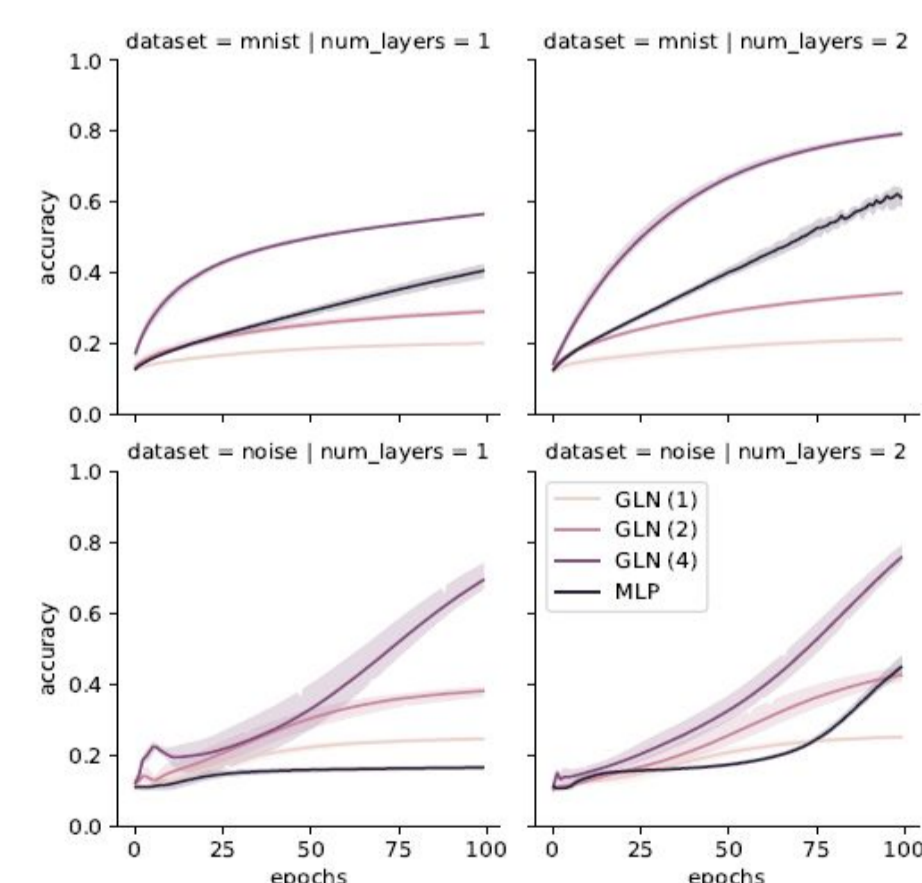


Figure 2. Output of a four layer network with random half-space contexts after training to convergence. Each box represents a non-bias neuron in the network, the function to fit is shown in black, and the output distribution learnt by each neuron is shown in colour (for example, red for the first layer and purple for the top-most neuron). All axes are identical, as labeled in the bottom left neuron. The dashed coloured lines represent the sampled hyperplane for each neuron.

- Can get a sense of model capacity by comparing ability to fit randomly shuffled or noisy labels.
- GLNs compare favourably with Deep ReLU networks.
- (Open Question) Can we characterize this formally?



## Linear Interpretability

$$\sigma\left(\underbrace{W_L(z)W_{L-1}(z)\dots W_1(z)}_{\text{multilinear polynomial of degree } L}\logit(p_0)\right),$$

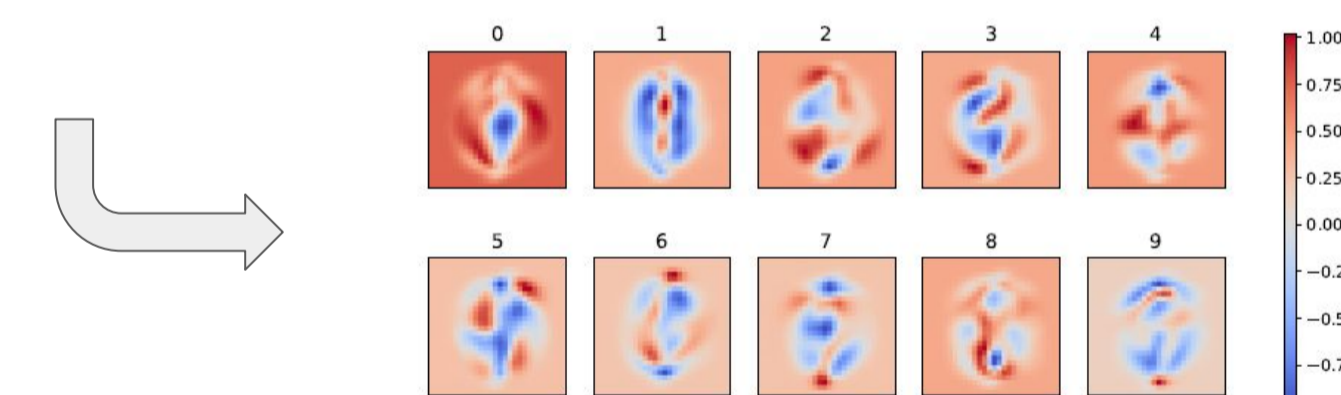


Figure 5. Saliency maps for constituent GLN binary classifiers of one-vs-all MNIST classifier after a single training epoch.

## Resilience to Catastrophic Forgetting

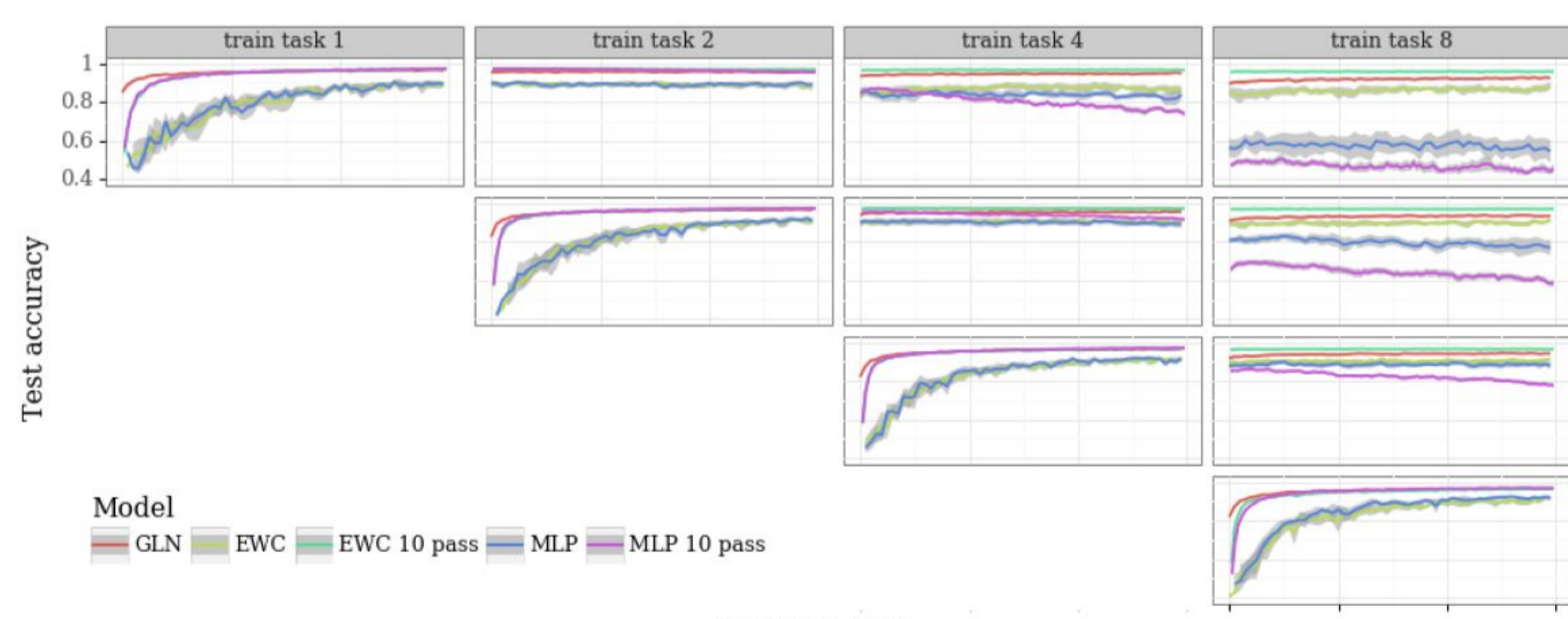


Figure 6. Retention results for permuted MNIST. Models are trained sequentially on 8 tasks (rows) and evaluated on all previously encountered tasks (columns). For example, the top-right plot indicates performance on Task 1 after being trained sequentially on Tasks 1 to 8 inclusive (not all tasks shown). Each model only trains for one epoch per task, with the exception of "EWC 10 pass" and "MLP 10 pass" (shrunk 10-fold on x axis). Error bars denote 95% confidence levels over 10 random seeds.

## Inductive Bias and Results



Figure 4. The effect of a single noisy XOR update (circled) on the decision boundaries of a halfspace gated GLN. Sampled hyperplanes for each gate are shown in white.

Inputs close in terms of cosine similarity will map to similar products of weight matrices!

- Single pass classification performance of GLNs matches general purpose batch techniques like SVMs, XGBoost, Deep ReLU Networks on UCI datasets.
- 98% accuracy with a single online pass over MNIST
- Matches SOTA NATs-per-image for autoregressive MNIST density modelling using just 1 pass!

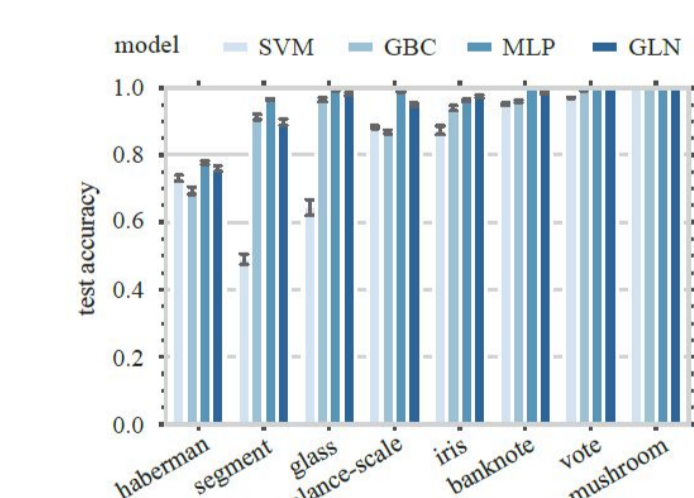


Figure 7. Online (single-pass) GLN classification accuracy on a selection UCI datasets, compared to three contemporary batch methods: Support Vector Machine, Gradient Boosting for Classification, Multi-Layer Perceptron trained for 100 epochs.

## Algorithm and Code

- Of course, the best way to build understanding and intuition is to see them in action: [github.com/deepmind/deepmind-research/tree/master/gated\\_linear\\_networks](https://github.com/deepmind/deepmind-research/tree/master/gated_linear_networks)

\* Equal contributions. Contact author: [aixi@google.com](mailto:aixi@google.com)

**Algorithm 1** GLN(Θ, z, p, x, y): update.  
Perform a forward pass and optionally update weights.

- 1: Input: GLN weights Θ = {σ<sub>ij</sub>}
- 2: Input: side info z, base predictions p ∈ [0, 1]<sup>K<sub>0</sub>-1</sup>
- 3: Input: binary target x, learning rate η ∈ (0, 1)
- 4: Input: boolean update (controls if we learn or not)
- 5: Output: estimate of P[x = 1 | z, β]
- 6: p<sub>0</sub> ← (Σ<sub>j=0</sub><sup>K<sub>0</sub>-1</sup> p<sub>j</sub>σ<sub>0j</sub>)
- 7: for i ∈ {1, ..., L} do {loops over layers}
- 8: p<sub>0i</sub> ← β
- 9: for j ∈ {1, ..., K<sub>i</sub>} do {loops over neurons}
- 10: p<sub>ij</sub> ← CLIP<sub>[0,1]</sub><sup>+</sup>{σ(w<sub>ij</sub>(z) + σ<sup>-1</sup>(p<sub>0i-1</sub>))}
- 11: if update then
- 12: Δ<sub>ij</sub> ← -η(p<sub>ij</sub> - x)σ<sup>-1</sup>(p<sub>0i-1</sub>)
- 13: w<sub>ij</sub>(z) ← CLIP<sub>[0,1]</sub><sup>+</sup>{w<sub>ij</sub>(z) + Δ<sub>ij</sub>}
- 14: end if
- 15: end for
- 16: end for
- 17: return p<sub>L1</sub>

**Algorithm 1 G-GLN:** inference with optional update

- 1: Input: base model / features {μ<sub>0j</sub>, σ<sub>0j</sub>}<sup>K<sub>0</sub>-1</sup>
- 2: Input: side information z ∈ Z, target y ∈ ℝ
- 3: Input: G-GLN weights {W<sub>ik</sub>}, learning rate η ∈ (0, 1)
- 4: Output: Gaussian PDF
- 5: for i ∈ {1, ..., L} do
- 6: for k ∈ {1, ..., K<sub>i</sub>} do
- 7: {w<sub>ik</sub>(z), σ<sub>ik</sub>(z)} ← W<sub>ik</sub>(z)
- 8: σ<sub>ik</sub><sup>2</sup> ← [Σ<sub>j=0</sub><sup>K<sub>i-1</sub></sup> w<sub>ij</sub>σ<sub>ij</sub><sup>2</sup>σ<sub>ik</sub><sup>-2</sup>]<sup>-1</sup>
- 9: μ<sub>ik</sub> ← σ<sub>ik</sub><sup>-1</sup>[Σ<sub>j=0</sub><sup>K<sub>i-1</sub></sup> w<sub>ij</sub>μ<sub>ij</sub>σ<sub>ij</sub>σ<sub>ik</sub><sup>-1</sup>]
- 10: W<sub>ik</sub>(z) ← PROJ<sub>[0,1]</sub>{W<sub>ik</sub>(z) - η∇<sub>W<sub>ik</sub></sub>ℓ(z)} // (if learning)
- 11: end for
- 12: end for
- 13: return N(μ<sub>L1</sub>, σ<sub>L1</sub><sup>2</sup>)

## Additional GLN Resources

- [1] Budden, David, et al. "Gaussian Gated Linear Networks", NeurIPS, 2020.
- [2] Wang, Jianan, et al. "A Combinatorial Perspective on Transfer Learning.", NeurIPS, 2020.
- [3] Sezener, Eren, et al. "Online Learning in Contextual Bandits using Gated Linear Networks", NeurIPS, 2020.
- [4] Sezener, Eren, et al. "Gated Linear Networks and Extensions, NeurIPS Beyond Backpropagation Workshop, 2020.
- [5] Veness, et al. "Online learning with gated linear networks" arXiv preprint arXiv:1712.01897, 2017.