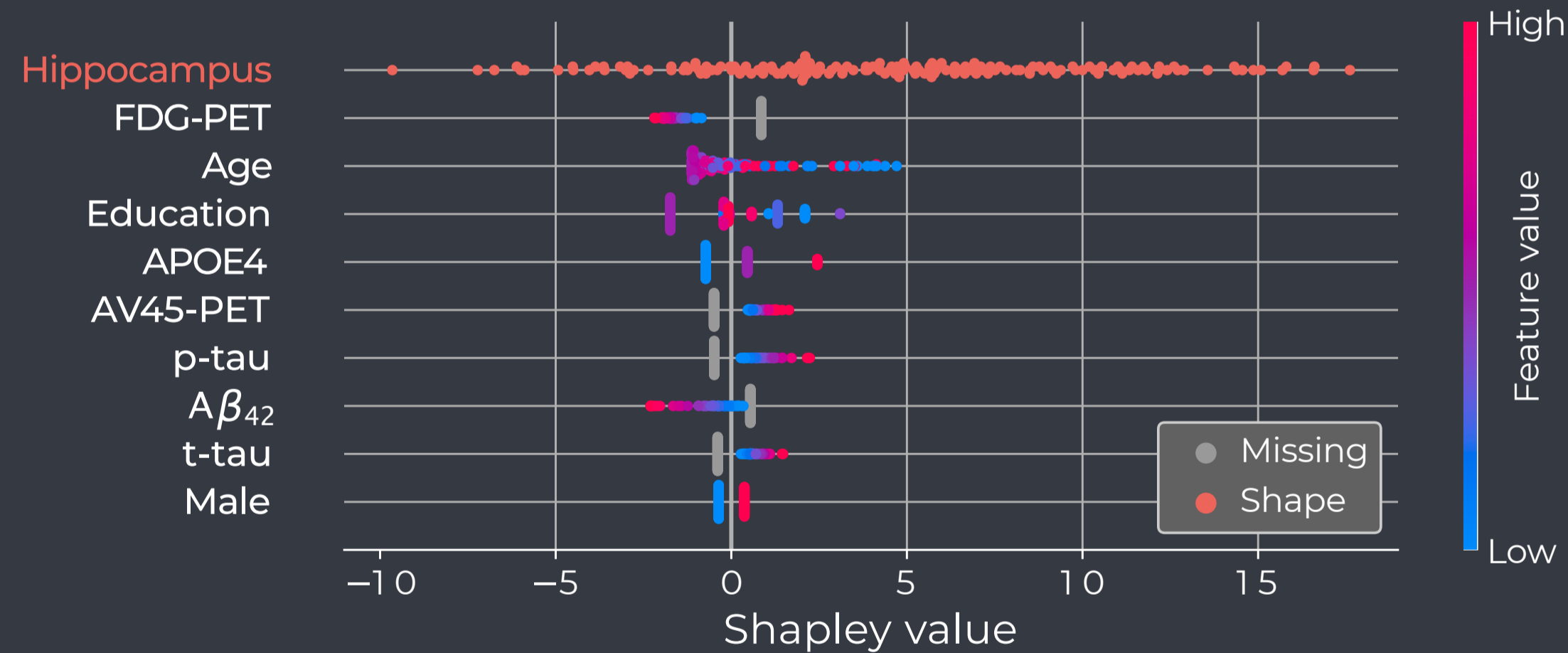
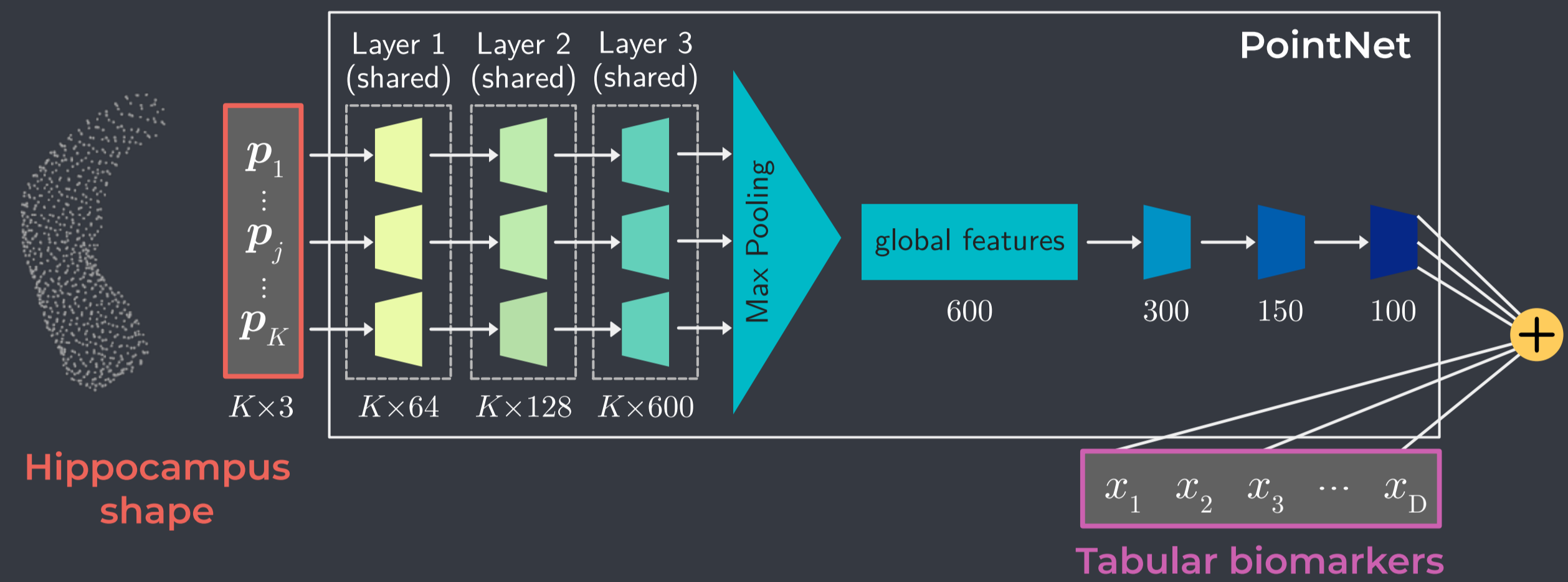


# Faithful and scalable **explanation of decisions** of a deep neural network from **anatomical shape** and **tabular biomarkers**.



## Scalable, Axiomatic Explanations of Deep Alzheimer's Diagnosis from Heterogeneous Data

Sebastian Pölsterl, Christine Aigner, and Christian Wachinger



## Explainable Artificial Intelligence (XAI)

- **Goal: Explain Alzheimer's diagnosis** made by a deep neural network  $f$  from **anatomical shape** and **tabular biomarkers**. Most existing work on XAI focus on CNN, but
  1. Inputs are *heterogeneous*,
  2. Shapes are *non-Euclidean*,
  3. Network *differs substantially from a standard CNN*.
- Axioms of explanations:

	Completeness	Null Player	Symmetry	Scale Invariance	Linear	Continuity	Implement. Invariance
Occlusion (Zeiler and Fergus, 2014)	x	✓	✓	✓	✓	x	✓
Guided Grad-CAM (Selvaraju et al., 2017)	x	✓	✓	✓	✓	x	✓
Layer-wise relevance prop. (Bach et al., 2015)	✓	✓	✓	✓	✓	✓	x
DeepLift (Shrikumar et al., 2017)	✓	✓	✓	✓	✓	✓	x
Integrated Gradients (Sundararajan et al., 2017)	✓	✓	✓	✓	✓	x	✓
Shapley Value (Shapley, 1953)	✓	✓	✓	✓	✓	✓	✓

## Methods

- Shapley Value of feature  $i$  of input  $\mathbf{z}$ :

$$\Delta_i = f(\mathbf{z}_{S \cup \{i\}}; \mathbf{z}_{\mathcal{F} \setminus S \cup \{i\}}^{\text{bl}}) - f(\mathbf{z}_S; \mathbf{z}_{\mathcal{F} \setminus S}^{\text{bl}})$$

$$s_i(\mathbf{z} | f) = \frac{1}{|\mathcal{F}|!} \sum_{S \subseteq \mathcal{F} \setminus \{i\}} |S|! \cdot (|\mathcal{F}| - |S| - 1)! \cdot \Delta_i$$

$$= \frac{1}{|\mathcal{F}|!} \sum_{k=0}^{|\mathcal{F}|-1} \sum_{\substack{S \subseteq \mathcal{F} \setminus \{i\} \\ |S|=k}} k! (|\mathcal{F}| - k - 1)! \cdot \Delta_i$$

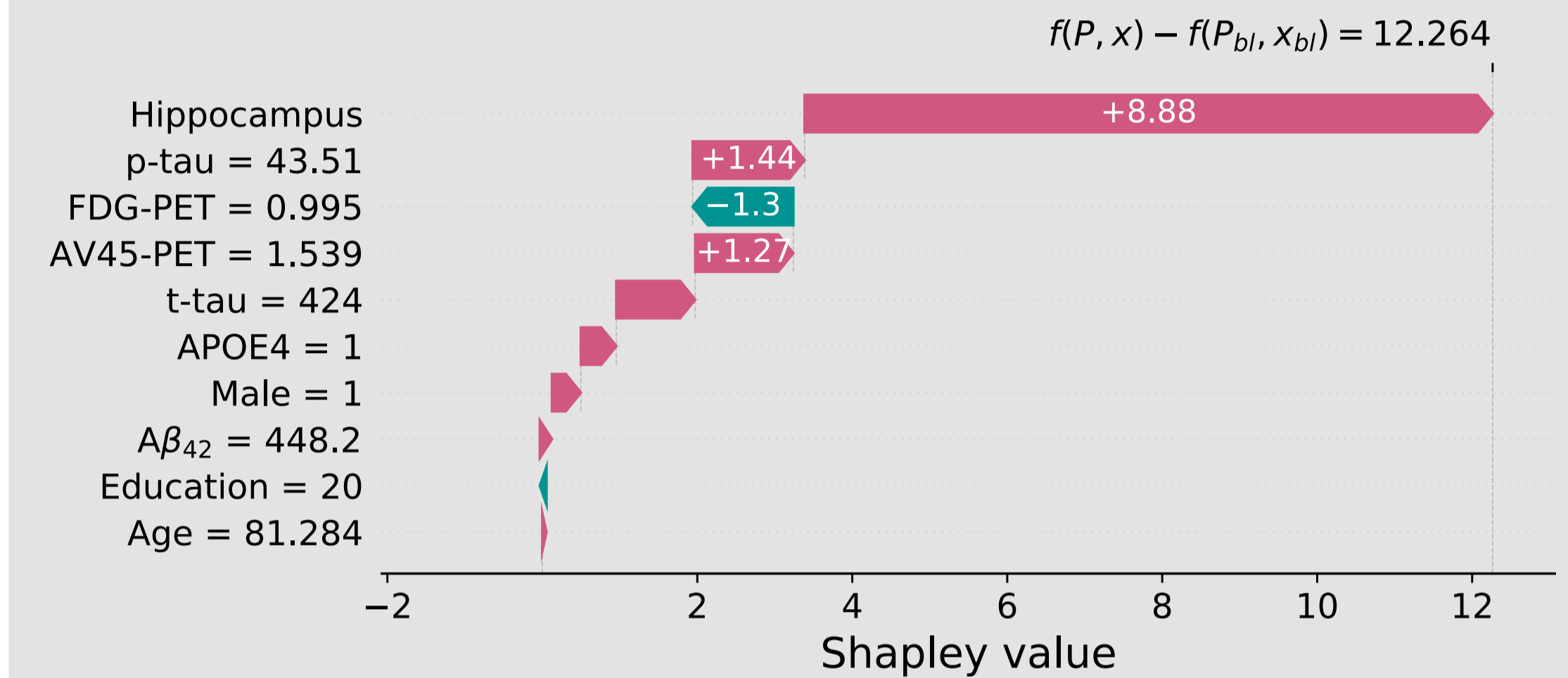
- ⊗ **Scales exponentially** in the number of features  $|\mathcal{F}|$ .
- $\mathbf{z}_{\mathcal{F} \setminus S}^{\text{bl}}$ : Simulate absence of point by replacing it with projection onto the convex hull.
- Approximate Shapley value:
  1. Represent output of first layer as *normal distribution* using sampling theory (Ancona et al., 2019; Cochran, 1977).
  2. Propagate distributions by converting layers into a *Lightweight Probabilistic Deep Network* (LDPN, Gast and Roth, 2018).

$$s_i(\mathbf{z} | f) \approx \underbrace{\mathbb{E}_k[f(\mathbf{z}_{S \cup \{i\}}; \mathbf{z}_{\mathcal{F} \setminus S \cup \{i\}}^{\text{bl}})]}_{\text{Output of LDPN}} - \underbrace{\mathbb{E}_k[f(\mathbf{z}_S; \mathbf{z}_{\mathcal{F} \setminus S}^{\text{bl}})]}_{\text{Output of LDPN}}$$

- ⊗ Approximate Shapley value **scales linearly**.

## Experiment: Real Data from ADNI

- Use Wide and Deep Network with left hippocampus shape ( $K = 1024$  points) and 9 tabular biomarkers.
- 1308 visits for training, 169 for hyper-parameter tuning.
- Balanced accuracy: 0.942 on test set (176 patients).
- Explanation of diagnosis for *individual* patient:



## Experiment: Synthetic Data

- Use PointNet to classify 100 point clouds of "X" and "I".
- Compare against *exact* Shapley value.

Method	MSE	SRC	NDCG	NE
SV Exact	0	1	1	65,538
Occlusion	16.1311	0.3180	0.8659	17
SV Sampling	0.0008	0.9505	0.9986	32,000
SV Sampling	0.0340	0.5440	0.9481	512
<i>Proposed</i>	0.0443	0.6918	0.9641	512

MSE: mean squared error. SRC: Spearman's rank correlation. NDCG: normalized discounted cumulative gain. NE: forward passes.