

Introduction

Invariant representation learning (IRL) [1]

expresses the inductive bias that a deep network's intermediate representations should also exhibit invariance to noise. Per sample x , we create noisy samples $\tilde{x}^1, \dots, \tilde{x}^K \sim \nu_x$ and modify the objective to penalize the distance d_l between their intermediate activations a_l :

$$\mathcal{L}_{\text{IRL}}(x, y) = \alpha \mathcal{L}(x, y) + \beta \mathcal{L}_{\text{noise}}(x, y) + \gamma \mathcal{L}_{\text{dist}}(x)$$

$$\mathcal{L}_{\text{noise}}(x, y) = \frac{1}{K} \sum_{k=1}^K \mathcal{L}(\tilde{x}^{(k)}, y),$$

$$\mathcal{L}_{\text{dist}}(x) = \sum_{\ell=1}^L \gamma_{\ell} d_{\ell}(a_{\ell}(x), \{a_{\ell}(\tilde{x}^{(k)})\}_{k=1}^K)$$

We take d_l to be a weighted sum of L2 and cosine distance (reduces degs. of freedom). We show IRL:

- can be interpreted from vicinal and structural risk minimization theory
- generalizes known stochastic and analytic regularizations for noise/adversarial robustness
- enables semi-supervised learning
- also applies to deep networks for computer vision and language modeling

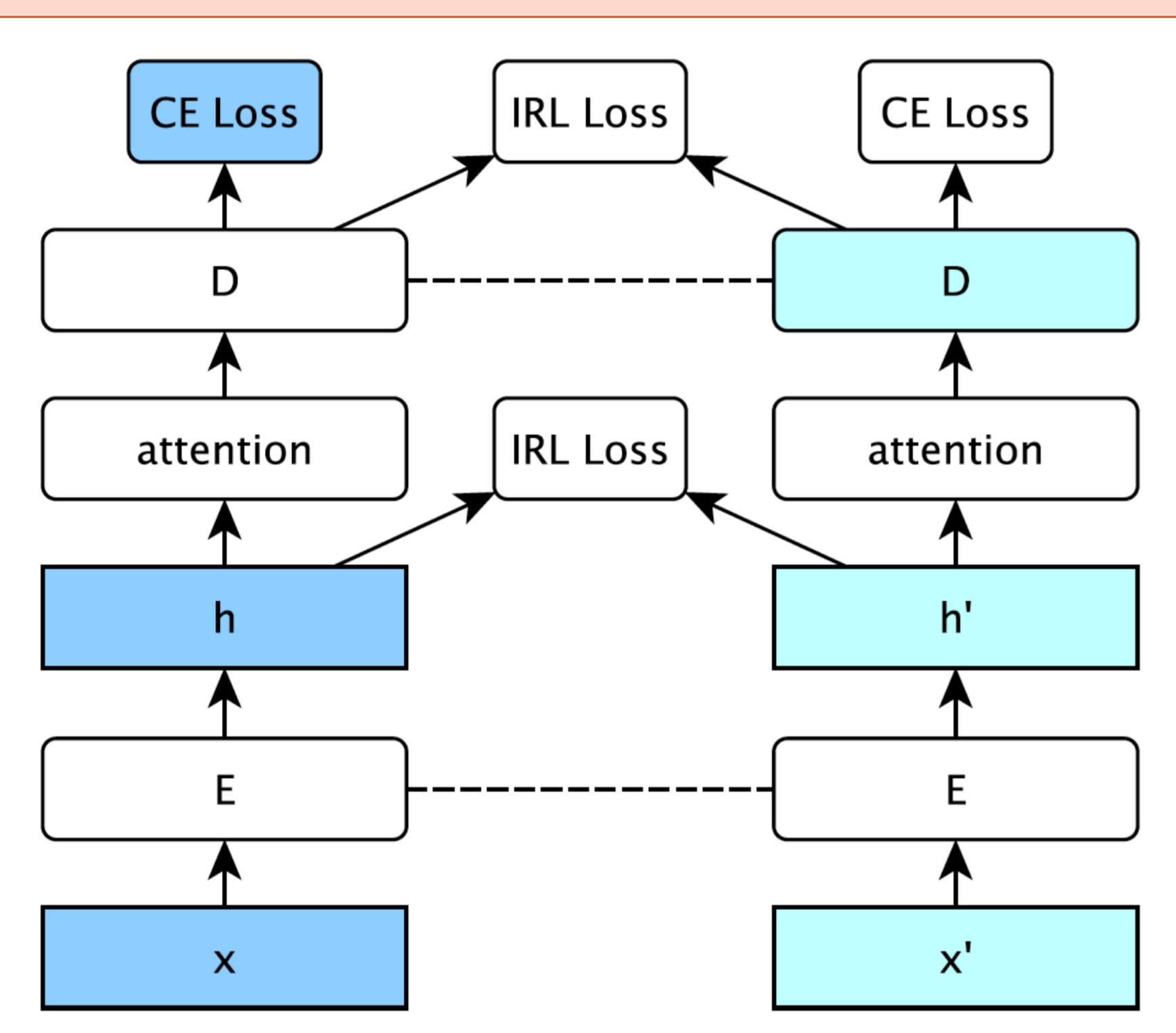


Figure 1. IRL loss with $K = 1$ for a sequence-to-sequence model.

Language Modeling

Our noise methods and baselines follow [9]; with probability $p = 0.2$, replace context words with a **unigram** draw or a special **blank** token. Here, IRL loss uniquely improves over the baseline where augmentation and IRL batching do not.

Computer Vision

We use IRL with Wide ResNet [8]. Their random crops, flips are applied before noising. We categorize noise into **in-domain** (brightness and contrast jitters, PCA noise), **out-of-domain** (hue and saturation jitters), **adversarial** (via FGSM [5]). IRL improves **5-8%** on the baseline trained in-domain, realized incrementally over data augmentation, IRL batching, and IRL loss. IRL also improves on regular and weighted adversarial training.

Table 2. Test set error on Wide ResNet-28-10 model. We apply IRL loss on the last four blocks.

(α, β, γ)	Method	Noise	CIFAR-10				CIFAR-100						
			none	in.	out.	$\epsilon=0.2$	$\epsilon=0.8$	none	in.	out.	$\epsilon=0.1$	$\epsilon=0.4$	
(1, 0, 0)	Baseline	std.	3.89	—	—	—	—	18.85	—	—	—	—	—
(1, 0, 0)	(ours)	std.	3.73	5.43	9.24	40.62	82.22	18.77	23.23	39.70	55.87	96.96	—
(0, 1, 0)	Data aug.	in.	3.79	4.67	9.48	42.58	86.19	18.68	20.33	38.69	57.57	96.88	—
(0.5, 0.5, 0)	IRL bat.	in.	3.66	4.55	9.50	42.94	86.98	18.30	20.18	38.54	57.60	97.57	—
(0.5, 0.5, 1)	IRL loss	in.	3.60	4.46	8.90	48.43	87.72	17.64	19.78	38.89	57.70	97.65	—
(0.5, 0.5, 0)	Adv. trn.	$\epsilon=0.4$	4.92	7.74	9.42	4.79	69.51	$\epsilon=0.2$	23.83	31.48	47.51	20.90	83.46
(0.8, 0.2, 0)	IRL bat.	$\epsilon=0.4$	5.14	8.70	11.45	5.04	71.61	$\epsilon=0.2$	20.52	25.60	39.85	21.06	80.40
(0.8, 0.2, 1)	IRL loss	$\epsilon=0.4$	4.55	7.81	10.32	4.63	69.47	$\epsilon=0.2$	22.49	28.84	45.20	20.26	88.63

References

1. Liang, D., et al. (2018). Learning noise-invariant representations for robust speech recognition. IEEE SLT
2. Chapelle, O., et al. (2001). Vicinal risk minimization. NeurIPS.
3. Bishop, C. (1995). Training with noise is equivalent to Tikhonov regularization. Neural Computation.
4. Zheng, S., et al (2016). Improving the robustness of deep neural networks via stability training. CVPR.
5. Goodfellow, I., et al (2015). Explaining and harnessing adversarial examples. ICLR.
6. Roth, K., et al (2018). Adversarially robust training through structured gradient regularization. ArXiv.
7. Kannan, H., et al (2018). Adversarial logit pairing. ArXiv.
8. Zagoruyko, S. and Komodakis, N. (2016). Wide residual networks. BMVC
9. Xie, Z., et al. (2017). Data noising as smoothing in neural network language models. ICLR

Theory

Empirical risk minimization approximates P_{data} :

$$\hat{P}_{\text{data}} = \frac{1}{N} \sum_{i=1}^N (\delta_{x^{(i)}} \times \delta_{y^{(i)}}).$$

Optimization in deep learning approximates ERM:

$$R(\theta) \approx \int \mathcal{L}(x, y; \theta) d\hat{P}_{\text{data}}(x, y) \approx \int \mathcal{L}(x, y; \theta) dP_{\text{data}}(x, y) \approx \frac{1}{B} \sum_{b=1}^B \mathcal{L}(x^{(b)}, y^{(b)}; \theta)$$

Vicinal risk minimization [2]: The estimate is improved by linear interpolation of δ_x with a noise model ν_x :

$$\int \mathcal{L}(x, y; \theta) d[\alpha \delta_x + \beta \nu_x](x) d\delta_y(y) \approx \alpha \mathcal{L}(x, y; \theta) + \frac{\beta}{K} \sum_{k=1}^K \mathcal{L}(\tilde{x}^{(k)}, y; \theta)$$

To mitigate catastrophic forgetting due to batching, this mix is interpolated with the \hat{y} predicted by the model:

$$\int (\mathcal{L}(x, y; \theta) + \gamma \mathcal{L}(x, f(x; \theta); \theta)) d\nu_x(x) d\delta_y(y) \approx \dots + \frac{\gamma}{K} \sum_{k=1}^K \mathcal{L}_{\text{c.e.}}(f(\tilde{x}^{(k)}; \theta), f(x; \theta))$$

for which symmetrizing and supervising multiple layers corresponds to L_{dist} . This stochastic form generalizes stability training [4], adversarial training [5], and logit pairing [7], although at best they all supervise at the logit level.

Structural risk minimization: As in [3], assuming a Gaussian noise model $\tilde{x} = x + \xi$ allows for analytic approximation. IRL and cross-entropy give:

$$\mathcal{L}_{\text{IRL}}(x, y) \approx (\alpha + \beta) \mathcal{L}_{\text{c.e.}}(x, y) + \beta \Omega_R + \gamma \Lambda_R,$$

$$\Omega_R = -\frac{1}{2} \sum_{j=1}^J y_j (\nabla_x \log f_j(x))^\top \Sigma (\nabla_x \log f_j(x)),$$

$$\Lambda_R = \sum_{\ell=1}^L \gamma_{\ell} \sum_{j=1}^J (\nabla_x (a_{\ell})_j(x))^\top \Sigma (\nabla_x (a_{\ell})_j(x)),$$

where $\Sigma = \text{Cov}(\xi)$. This analytic form generalizes works like [3] which assume isotropy/diagonal Σ , or structural gradient regularization [6] which computes a running estimate.

Table 1. Perplexities of a two-layer, 1500 unit, word-level LSTM.

Method	Noise	PTB		WikiText-2	
		val.	test	val.	test
Baseline	none	81.6	77.5	—	—
(ours)	none	80.5	77.5	96.9	92.1
Data aug.	uni.	85.8	82.9	106.7	100.0
IRL bat.	uni.	80.2	77.4	97.5	92.6
IRL loss	uni.	77.2	74.4	99.2	93.7
Data aug.	blank	78.8	75.3	98.6	92.3
IRL bat.	blank	80.3	76.7	97.7	93.8
IRL loss	blank	75.1	71.8	94.0	88.6

Semi-Supervision

For unlabeled data, one can continue viewing ν_x as kernel density estimate of P_{data} . Furthermore, [2] notes the model's best estimate can be used to approximate the unknown y . This gives

$$\int \mathcal{L}(x, f(x; \theta); \theta) d\nu_x(x)$$

which corresponds to L_{dist} as before. Using the same CIFAR model and in-domain noise on 4,000 labeled gives **21.5%** and **60.7%** test error. Semi-supervised training with L_{dist} on unlabeled samples improves this to **18.2%** and **59.6%** respectively.

Activations and Gradients

In our CIFAR models, data augmentation already induces nearby intermediate representations, further improved by IRL. As predicted, per-layer Jacobians are reduced, especially in the last four layers relative to augmentation and IRL batching.

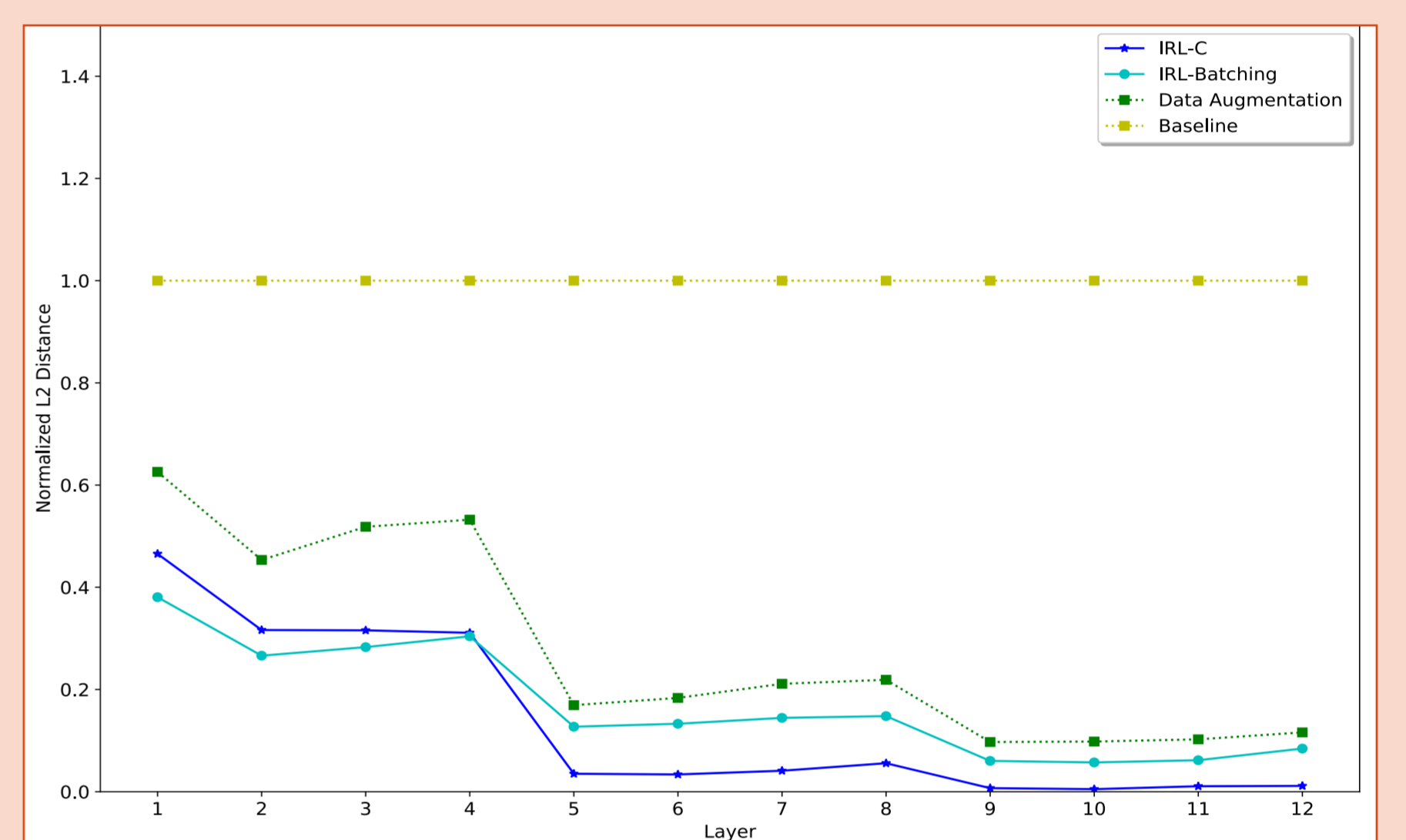


Figure 2. Normalized avg. L2 distance between x and x'

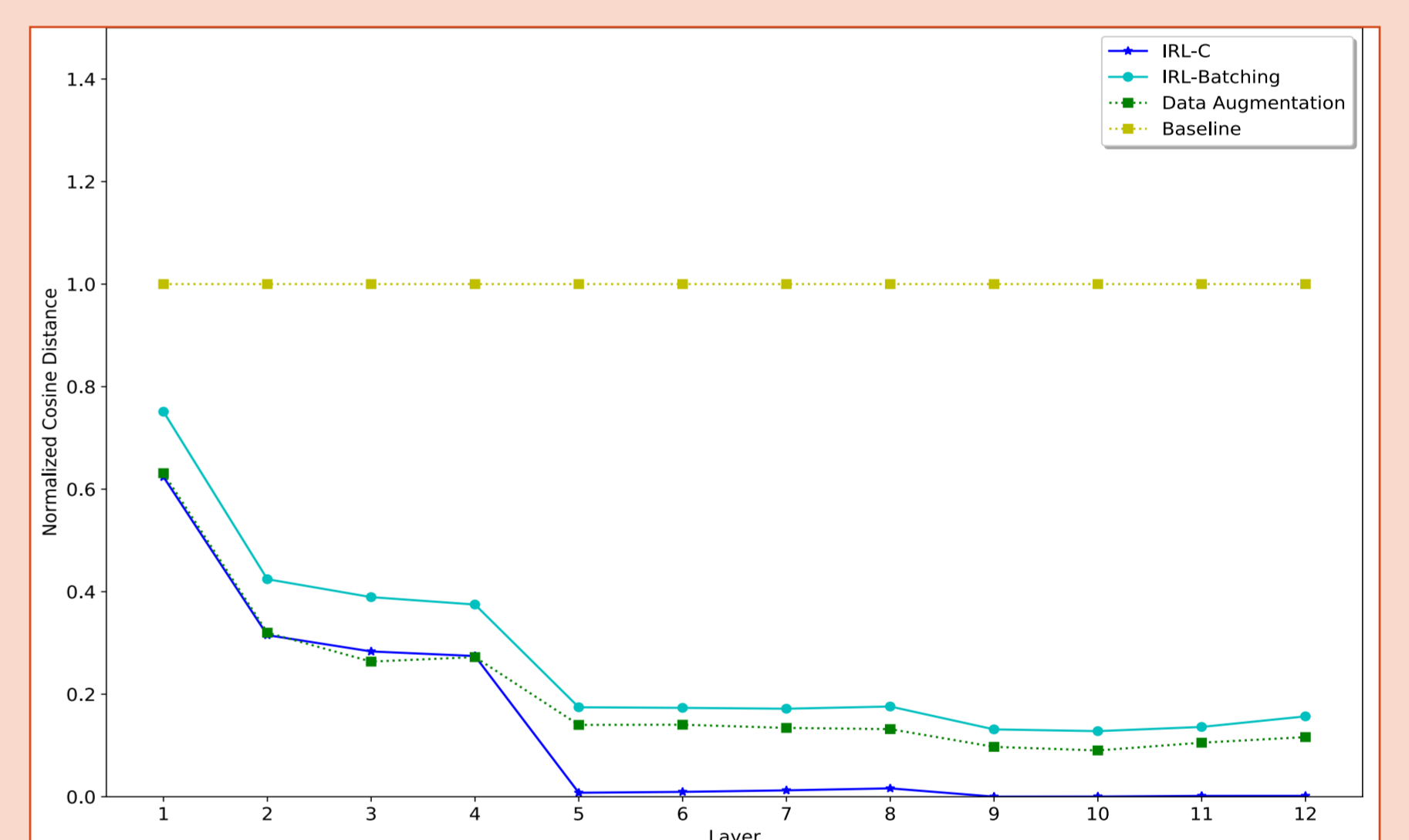


Figure 3. Normalized avg. cosine distance between x and x'

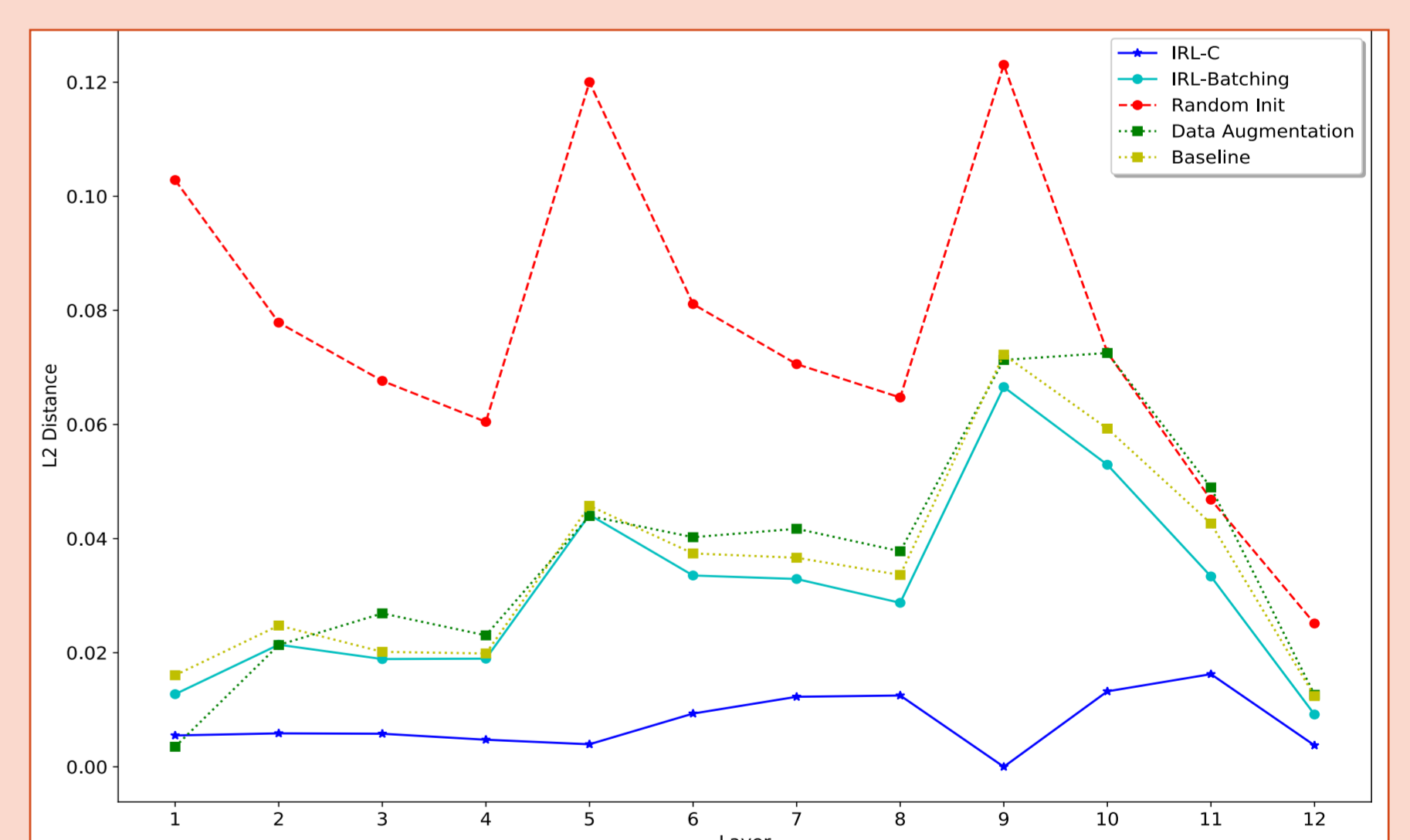


Figure 4. Jacobian norms per layer on the test set

Speech Recognition

Our noise model is additive from MUSAN, as in [1]. Here we use WSJ, for which IRL training improves accuracy in three of five out-of-domain conditions.

Table 3. Character errors of 4-4 Enc-Dec. on WSJ on noise

Method	WSJ (eval92)				
	none	RIR	speech	v.up	v.down tel.
Regular	11.2	67.5	130.0	33.8	28.0
IRL bat.	13.0	45.9	83.5	29.7	64.7
IRL loss	12.3	48.2	58.3	29.6	73.0