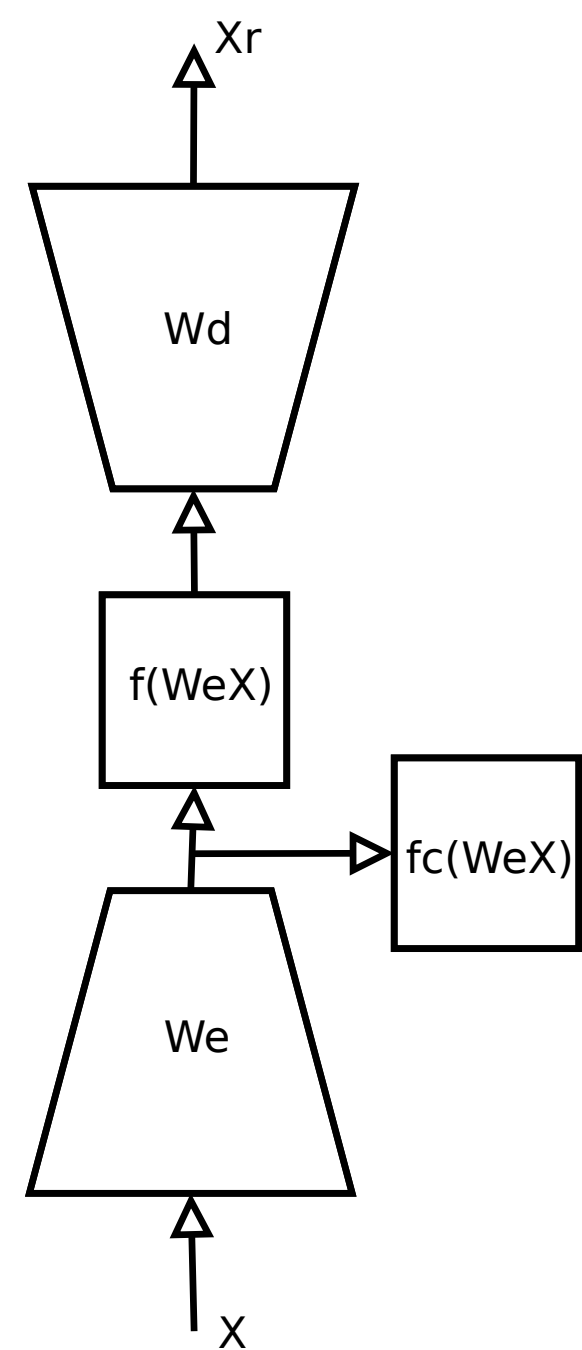# Saturating Auto-Encoders

## Rostislav Goroshin and Yann LeCun

Courant Institute of Mathematical Sciences, New York University

{goroshin,yann}@cims.nyu.edu

## Auto-Encoders Learn Manifolds

- Data lies on a low-dimensional manifold embedded in a higher dimensional ambient space
- Reconstruction objective ensures low reconstruction error in portions of the input space where data is densely distributed
- **Regularizers should raise reconstruction error for inputs not near the data manifold**
- Analogous to minimizing the partition function in max-likelihood models
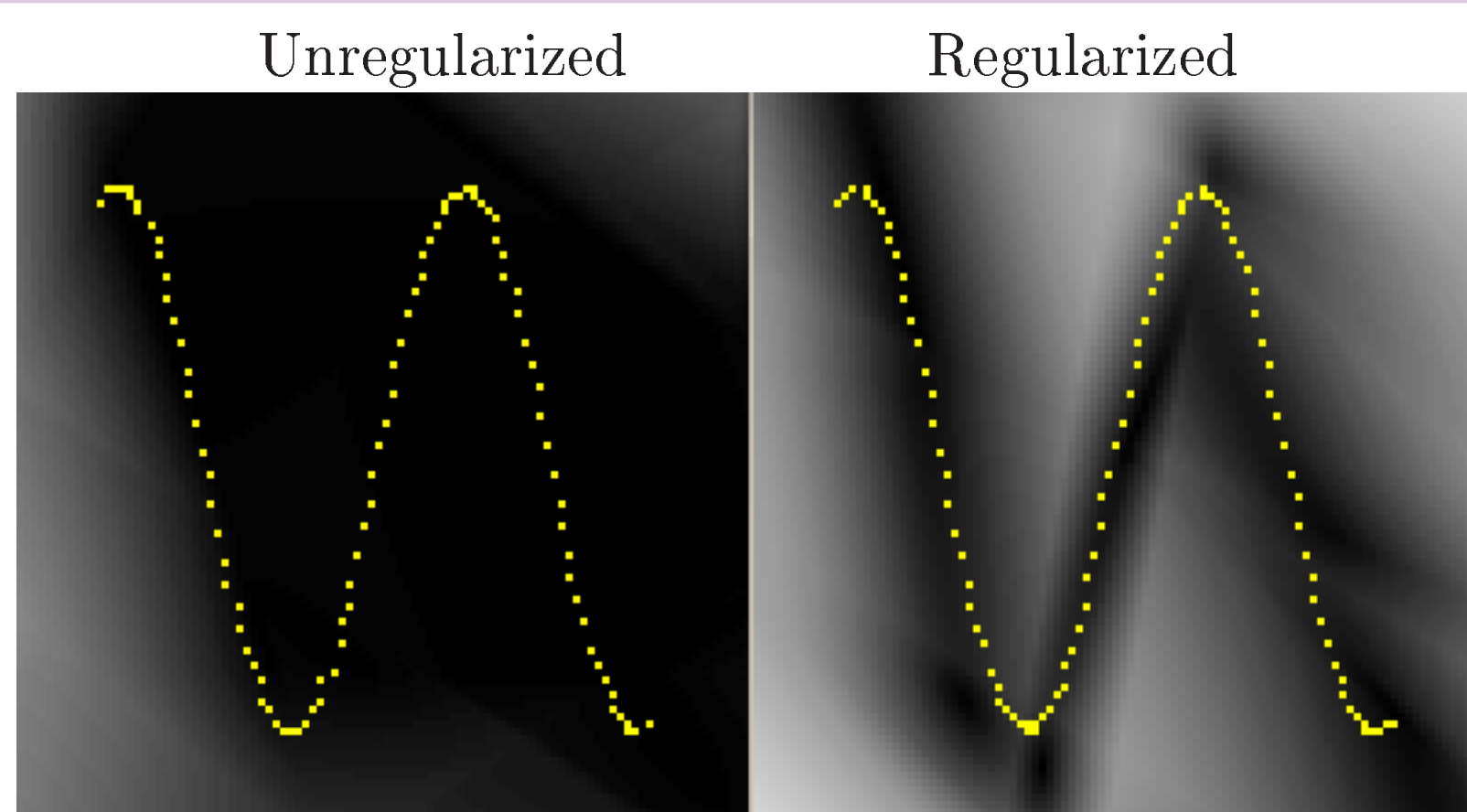


## Latent State Regularization

Latent state regularization is a method of introducing an information bottleneck, which is more intuitive than regularizing the weights directly.
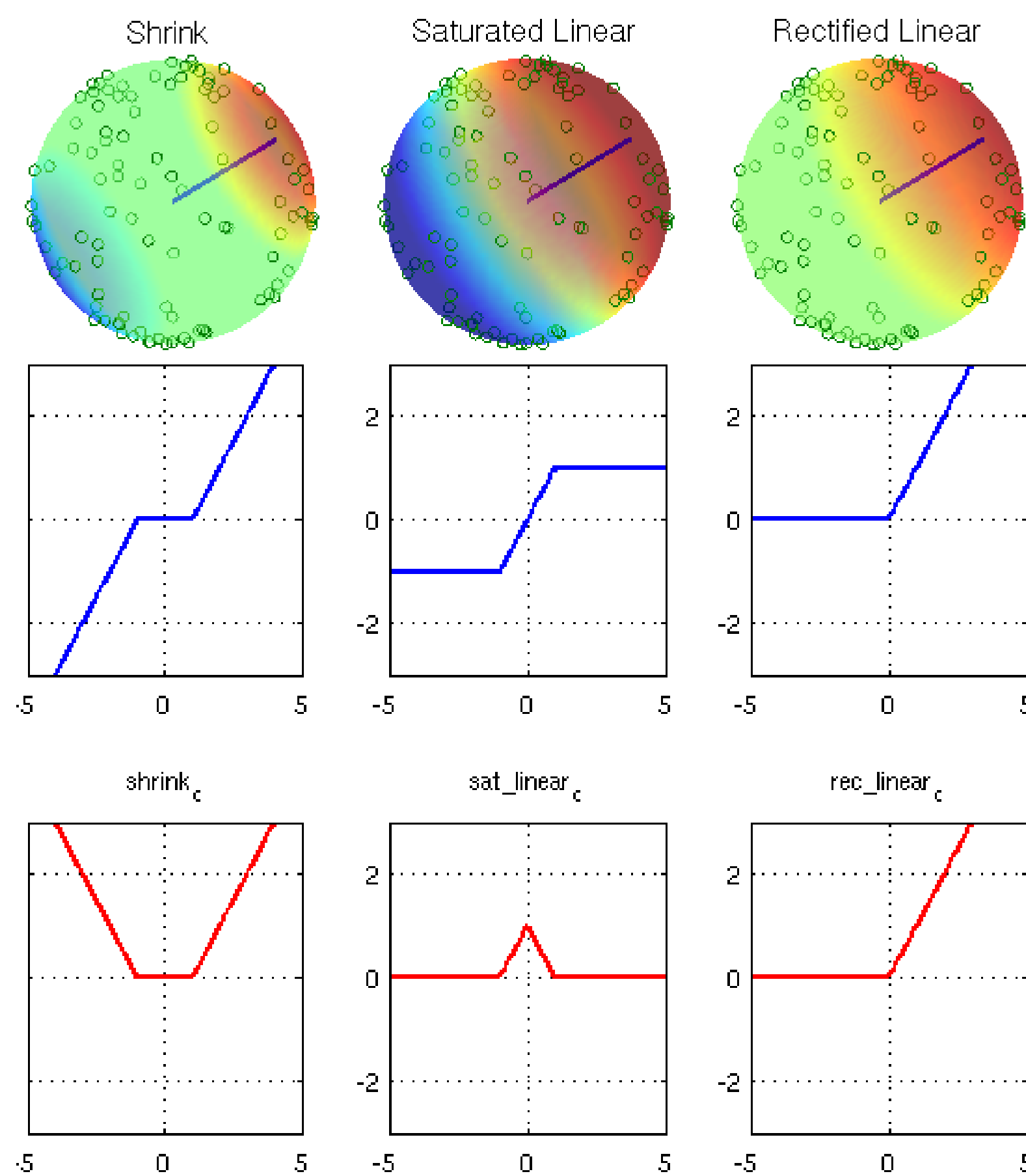
- Unsupervised Setting: Sparsity ($L_1$) and Contractive Regularization
- Supervised Setting: Penalize the Surface Area of the Decision Boundary

## Toy-Manifold Example

| Unregularized | Regularized |
|---|---|



## Regularization via Saturation

- Consider activation functions with flat (zero-gradient) regions
- These activation functions lose their ability to reconstruct the input when activations occur in the flat regions
- Reconstruction error will increase quadratically as we move away from the manifold
- We associate with each activation function a complimentary function which encourages activations in the flat regions



_Complimentary Function_
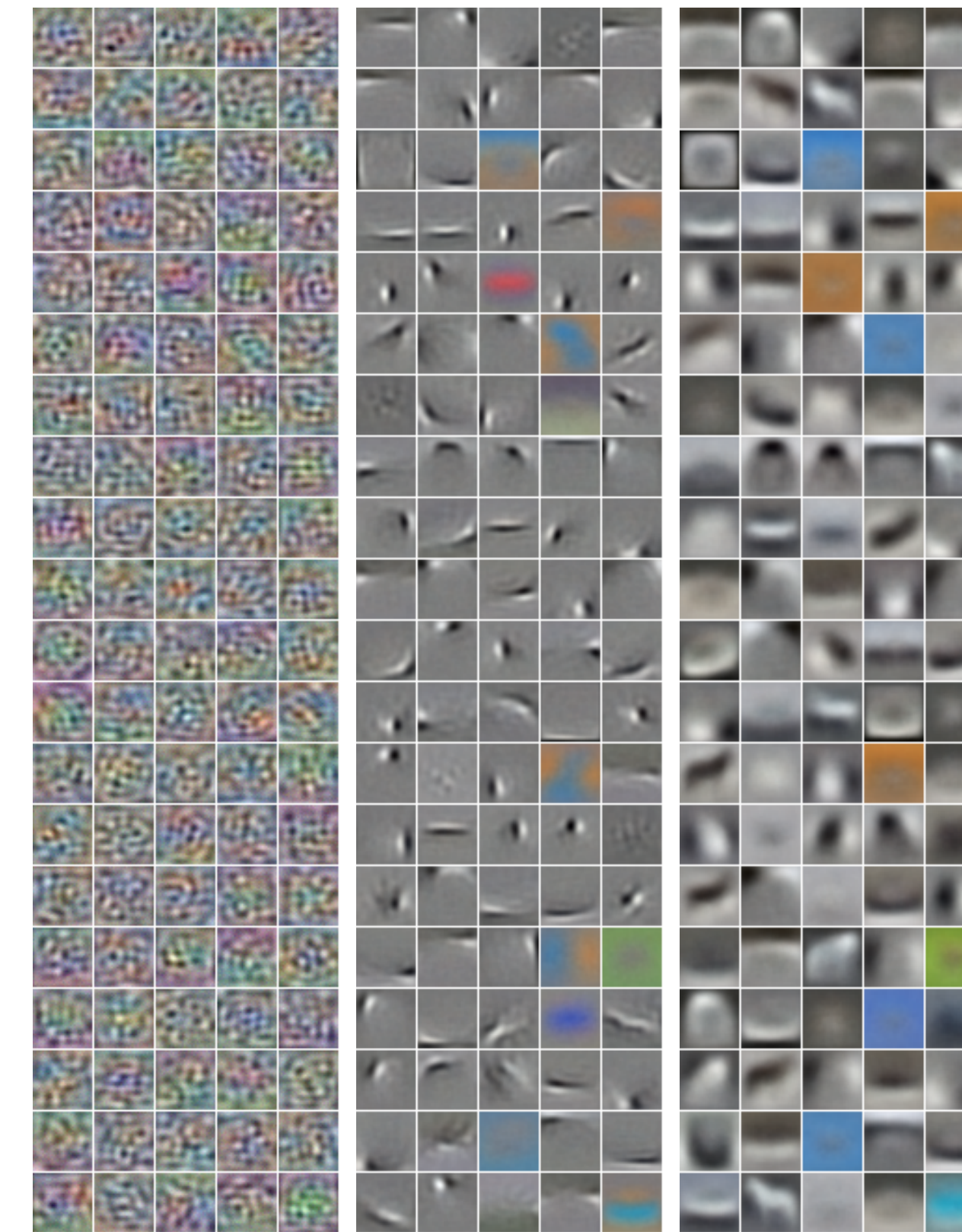
$$f_c(z) = \inf_{z' \in S} |z - z'|$$

_Loss Functional_

$$L = \sum_{x \in D} \frac{1}{2} \|x - (W^d f(W^e x + B^e) + B^d)\|^2$$

$$+ \alpha \sum_{i=1}^{d_h} f_c(W_i^e x + b_i^e)$$

## SATAE-shrink on CIFAR-10



## SATAE-sat-linear on CIFAR-10



## Progressively Increasing $\alpha$



## Relation to Other Regularizers

_Connection to Sparse-Auto-Encoders_

Note that $shrink_c(W^e x + b^e) = abs(shrink(W^e x + b^e))$. For $shrink$-nonlinearity, our regularizer corresponds to $L_1$ penalty on the activations.

_Connection to Contractive-Auto-Encoders_

$$\sum_{ij} \left( \frac{\partial h_i}{\partial x_j} \right)^2 = \sum_i^{d_h} \left( f'(\sum_{j=1}^d W_{ij}^e x_j + b_i)^2 \|W_i^e\|^2 \right)$$

The first term in the above equation tries to adjust the weights so as to push the activations into the low gradient (saturation) regime of the nonlinearity, but is only defined for differentiable activation functions. Therefore the CAE indirectly encourages operation in the saturation regime.

## Extension to $C^1$

We employ the concept of average variation over a finite interval to extend the definition of complimentary activation functions to differentiable functions.
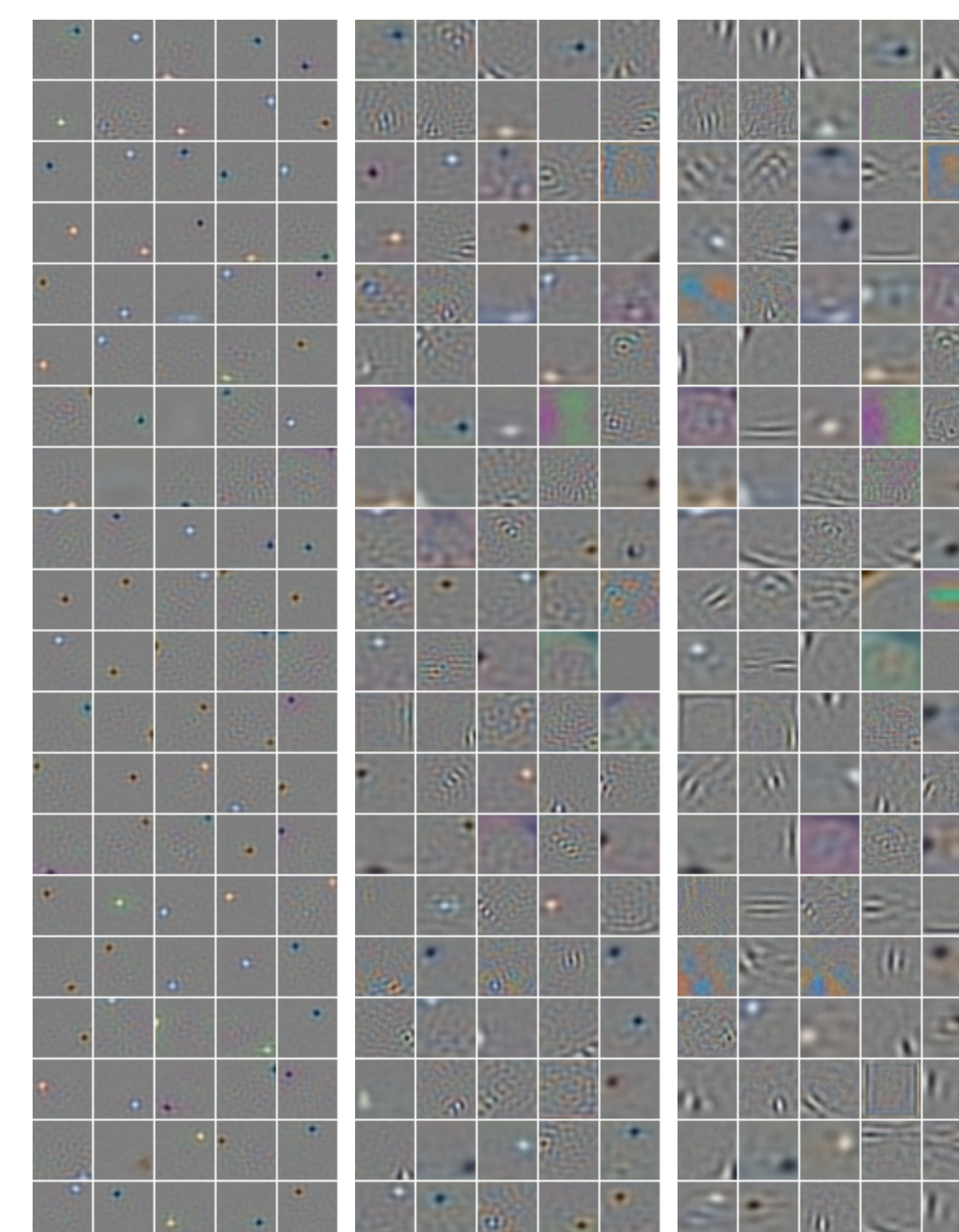
## References

1. Rifai, S. and Vincent, P. and Muller, X. and Glorot, X. and Bengio, Y. Contractive auto-encoders: Explicit invariance during feature extraction, _Proceedings of the Twenty-eight International Conference on Machine Learning, ICML 2011_
2. Karol Gregor and Yann LeCun: Learning Fast Approximations of Sparse Coding, Proc. _International Conference on Machine learning (ICML'10)_, 2010
3. Marc'Aurelio Ranzato, Christopher Poultney, Sumit Chopra and Yann LeCun. Efficient Learning of Sparse Representations with an Energy-Based Model, in J. Platt et al. (Eds),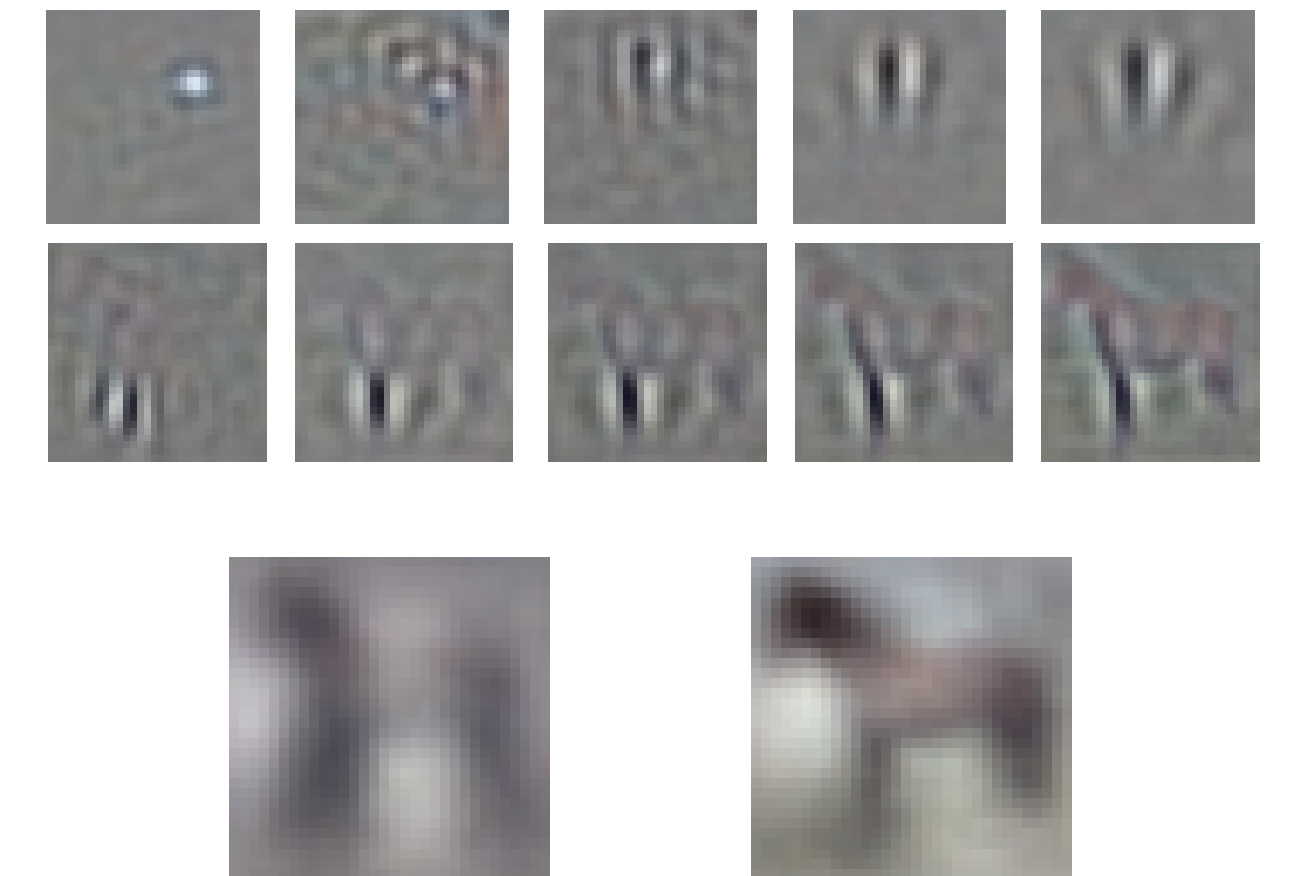 _(NIPS 2006)_, 19, MIT Press, 2006.