

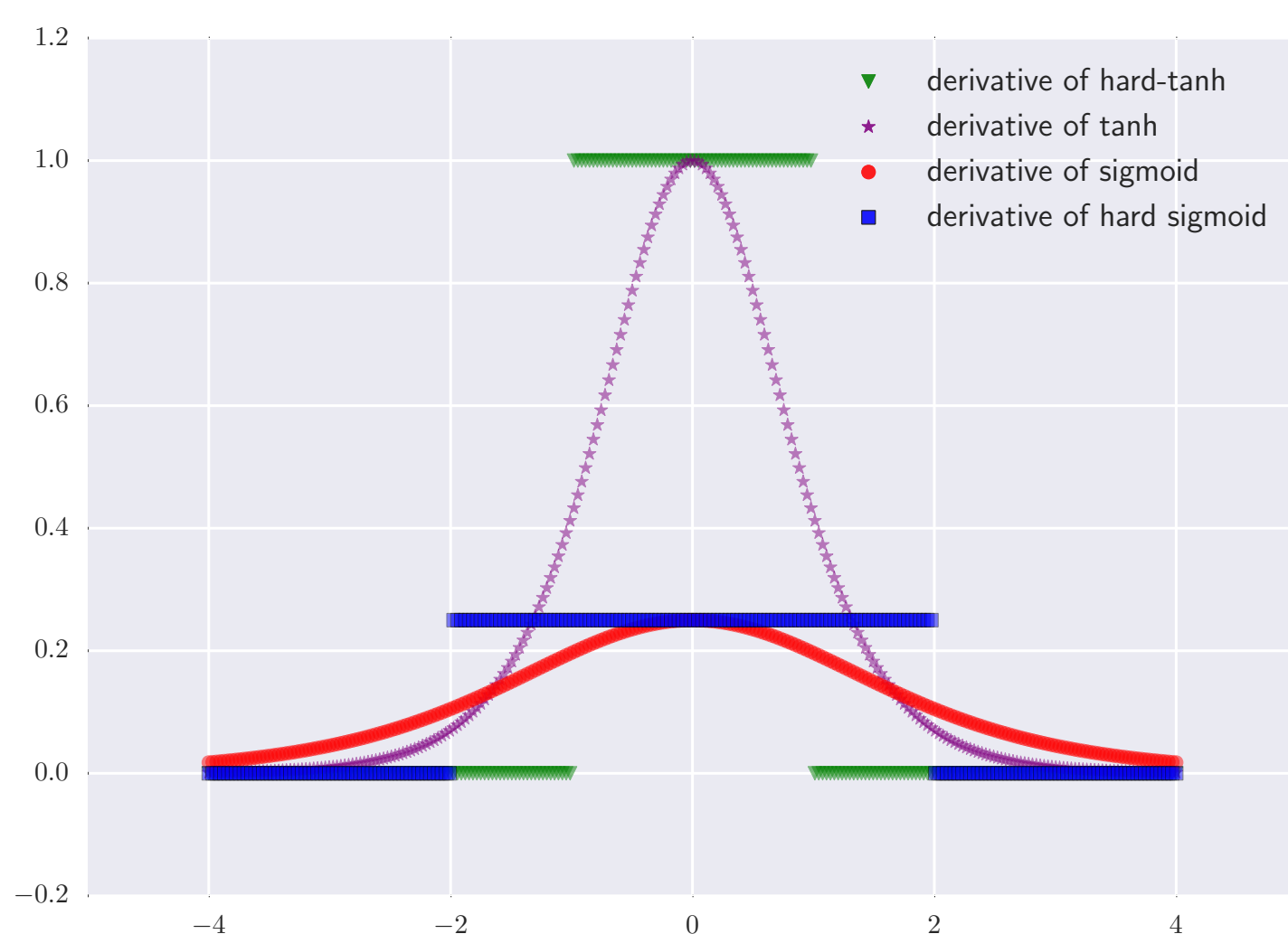
## MOTIVATION

- Common nonlinear activation functions can have training difficulties.
- Logistic functions (sigmoid and tanh) can be difficult to train.
- Piecewise linear activation functions (i.e. ReLU) are easier to optimize.

## OUR CONTRIBUTIONS

- Applying piecewise linear activation functions to gates of the recurrent models, i.e. LSTMs.
- Investigation of injecting noise to the activations.
- An efficient way to learn the std of noise for each unit.
- Annealing the activation noise can have a continuation effect.

## SATURATING ACTIVATIONS



**Definition 1.** *Soft-saturating activation:* An activation function softly saturates if it converges to a particular value as  $x \rightarrow \infty$  and/or  $x \rightarrow -\infty$ .

**Definition 2.** *Hard-saturating activation:* An activation function hardly saturates if it becomes constant when its input gets larger than a threshold  $c$ .

Linearize the activation function and clip it at the threshold:

$$\text{hard-sigmoid}(x) = \max(\min(u^s(x), 1), 0)$$

$$\text{hard-tanh}(x) = \max(\min(u^t(x), 1), -1)$$

## NOISY ACTIVATIONS (UNBIASED)

$h(x)$  : hard activation function.

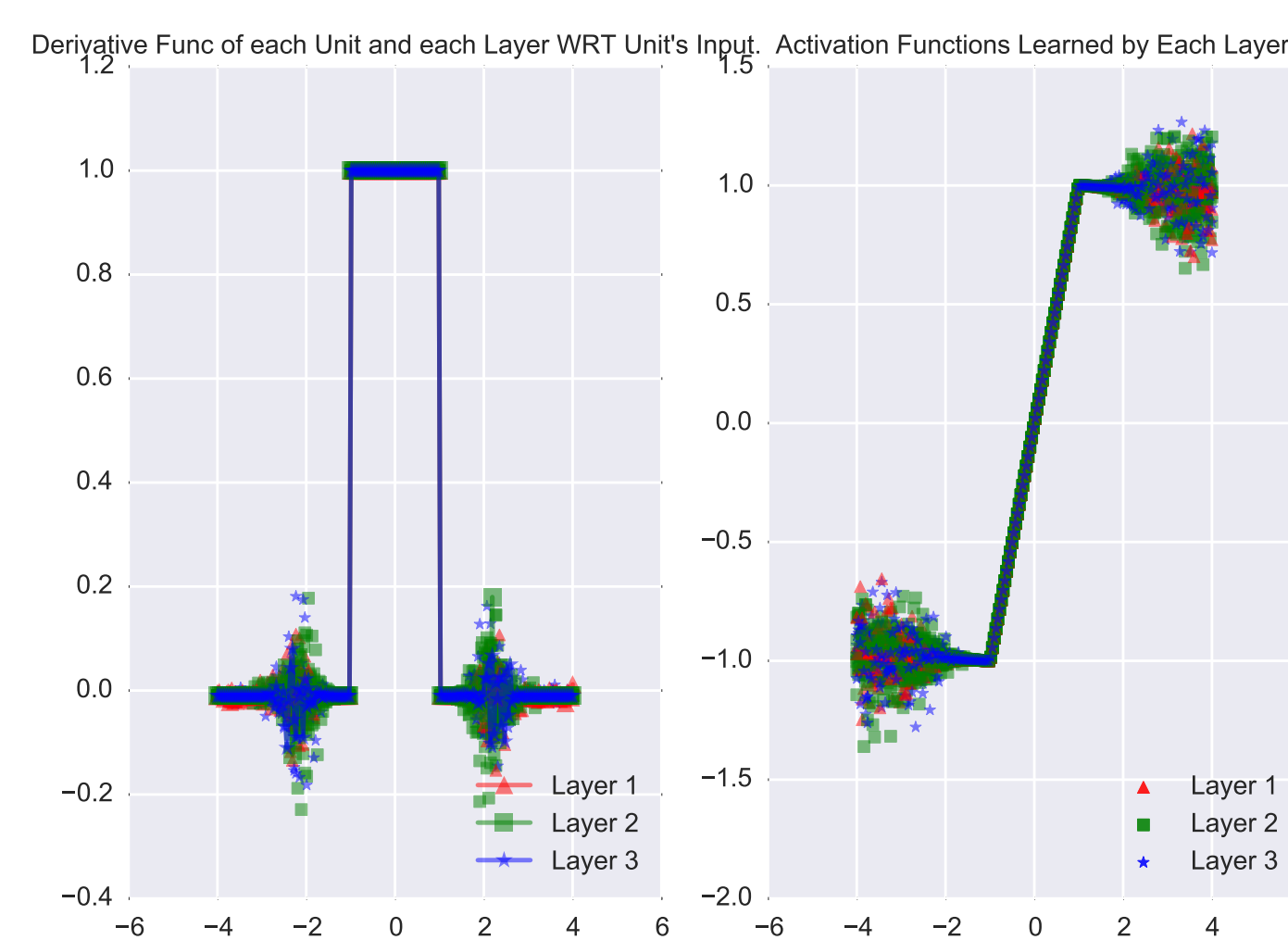
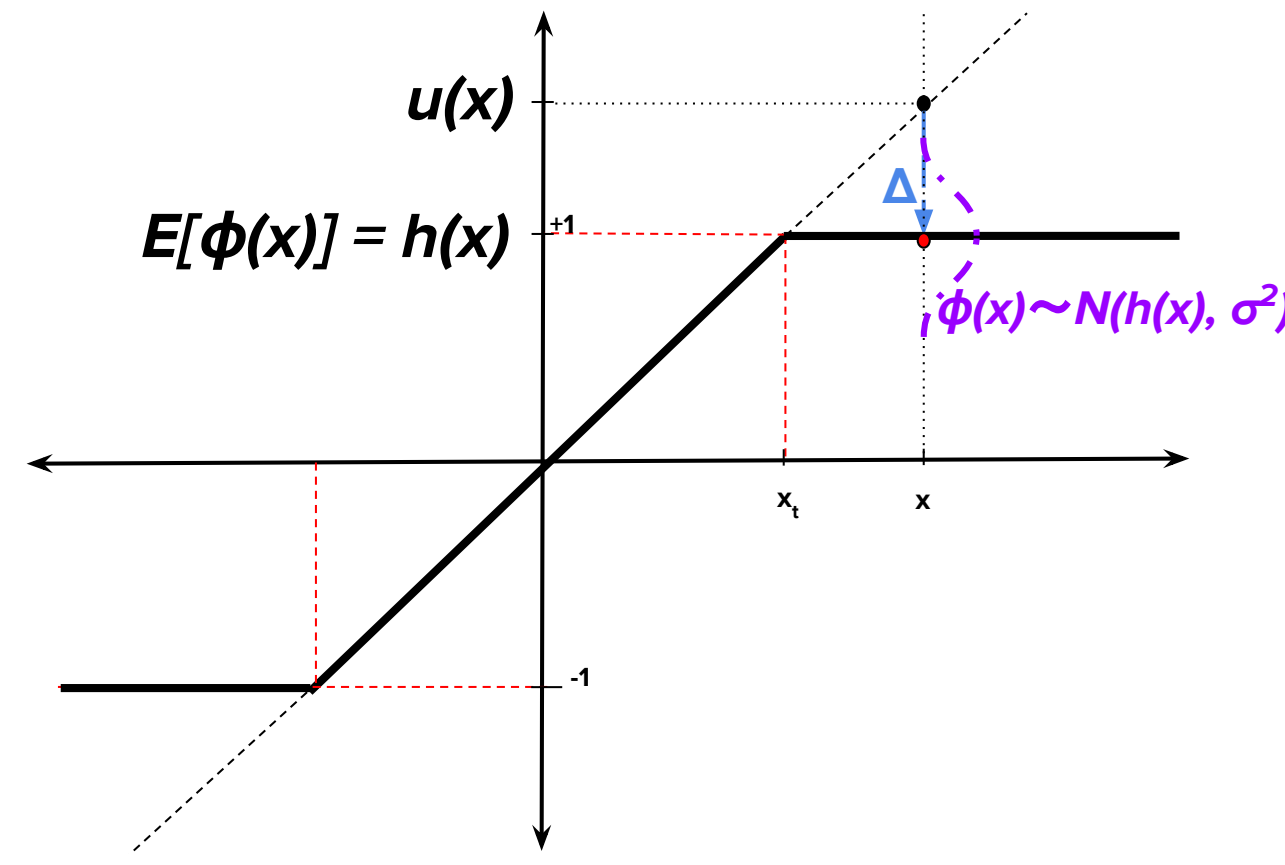
$u(x)$  : soft activation function.

$$\phi(x, \xi) = u(x) + s$$

$$s = \mu + \sigma \xi$$

$$E_{\xi \sim \mathcal{N}(0, 1)} \approx h(x)$$

## NOISY ACTIVATIONS (BIASED)



### Injecting Biased Noise:

$$d(x) = -\text{sgn}(x) \text{sgn}(1 - \alpha)$$

For  $\epsilon = |\xi|$ ,

$$s = \mu(x) + d(x) \sigma(x) \epsilon,$$

$$\phi(x, \xi) = \alpha h(x) + (1 - \alpha) u(x) + d(x) \sigma(x) \epsilon.$$

Use the expectation of the noise at the test time:

$$E[\phi(x, \xi)] = \alpha h(x) + (1 - \alpha) u(x) + d(x) \sigma(x) E[\epsilon].$$

### Injecting Noise at the Input:

Noise injection to the input of the activation function can be written as:

$$\phi(x, \xi) = h(x + s).$$

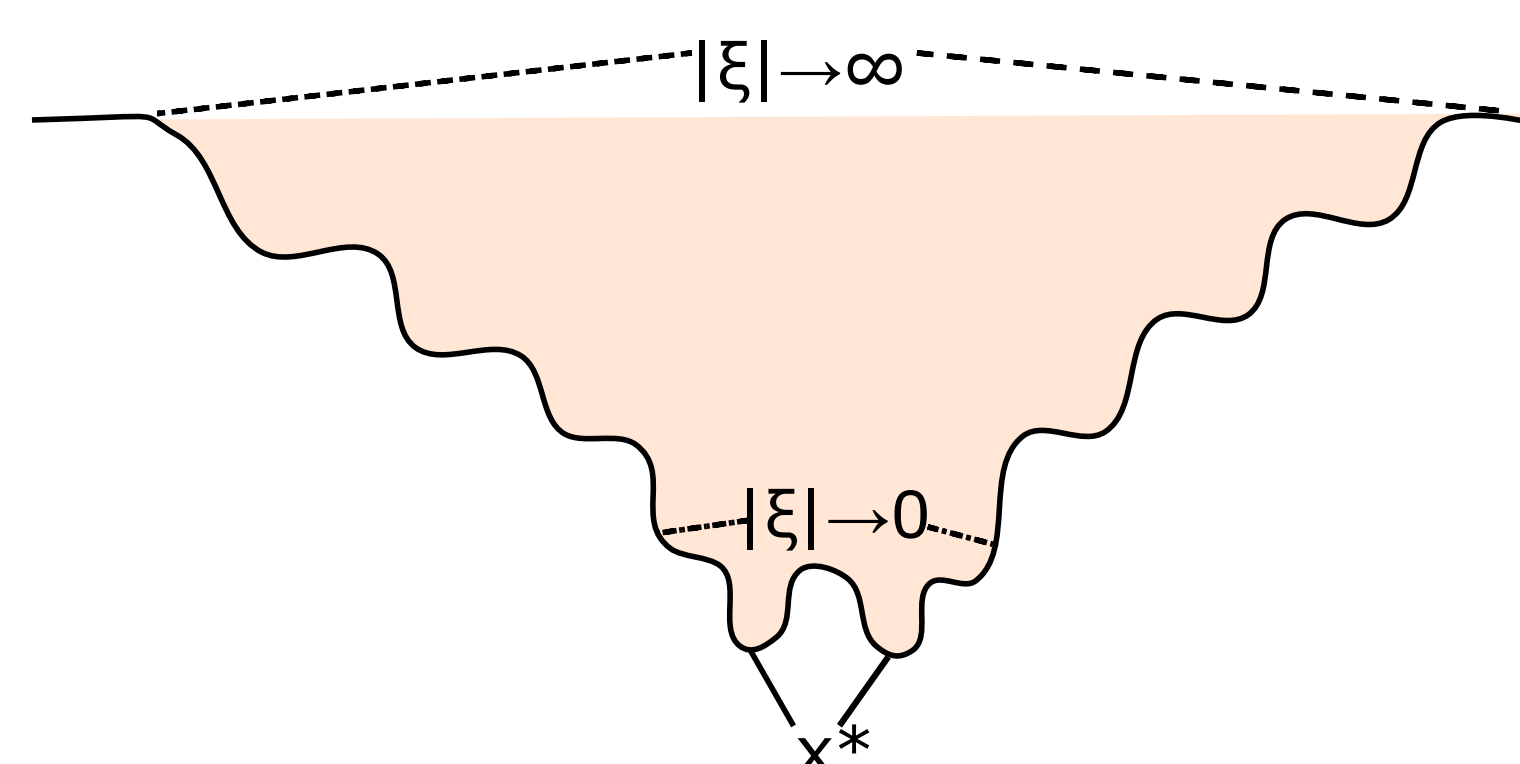
and  $s$  can have different formulations, e.g.  $s = \sigma \xi$ ,  $s = \sigma(x) \xi$  or  $s = \mathbf{1}_{|x| \geq |x_t|} (\sigma \xi)$

## ANNEALING THE NOISE

Start with large noise resulting in larger exploration and anneal the noise:

$$\lim_{|\xi| \rightarrow \infty} \left| \frac{\partial \phi(x, \xi)}{\partial x} \right| \rightarrow \infty.$$

A pathological one-dimensional case for SGD:



## EXPERIMENTS

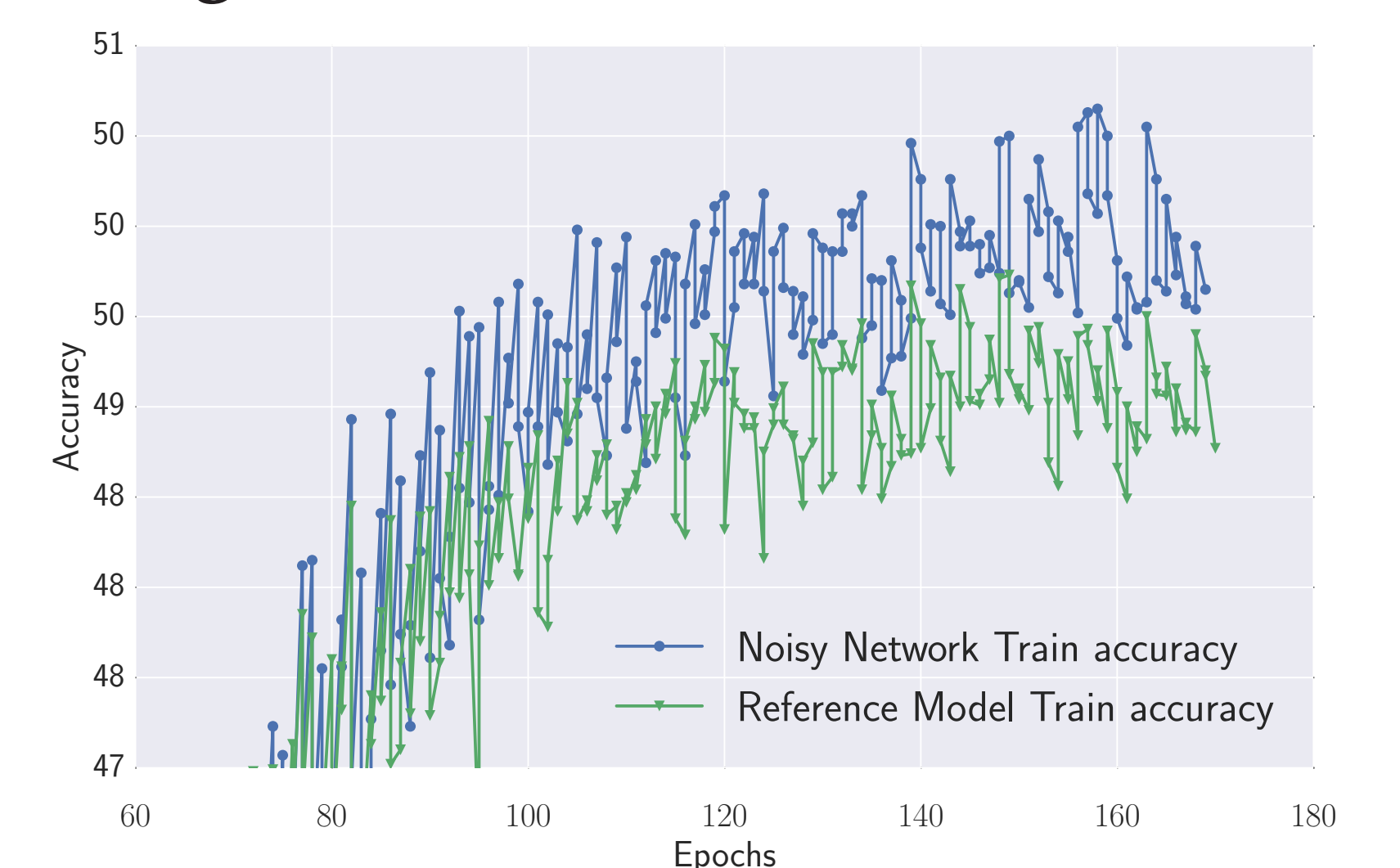
Used the same hyperparameters with base-lines.

- NAN** - Normal Noise at the output.
- NAH** - Half-Normal (biased) Noise at the output.
- NANI** - Normal Noise at the Input.
- NANIL** - Normal Noise with Learned  $\sigma(x)$  at the Input.
- NANIS** - Normal Noise at the Input when unit Saturates.

### Neural Machine Translation

	Valid nll	BLEU
Sigmoid and Tanh NMT (Reference)	65.26	20.18
Hard-Tanh and Hard-Sigmoid NMT	64.27	21.59
Noisy (NAH) Tanh and Sigmoid NMT	<b>63.46</b>	<b>22.57</b>

### Learning to Execute



### Image Caption Generation

\*Image caption generation is both without dropout.

	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	Test NLL
Soft (Ref.)	67	44.8	29.9	19.5	18.9	40.33
Soft (NAH)	66	<b>45.8</b>	<b>30.69</b>	<b>20.9</b>	<b>20.5</b>	40.17
Soft (NAH*)	64.9	44.2	<b>30.7</b>	<b>20.9</b>	20.3	<b>39.8</b>
Soft (NANI)	66	45.0	30.6	20.7	<b>20.5</b>	40.0
Soft (NANIL)	66	44.6	30.1	20.0	<b>20.5</b>	39.9
Hard	67	45.7	31.4	21.3	19.5	-

### PennTreeBank Experiments

	Valid ppl	Test ppl
Noisy LSTM + <b>NAN</b>	111.7	<b>108.0</b>
Noisy LSTM + <b>NAH</b>	112.6	<b>108.7</b>
LSTM (Reference)	119.4	115.6

### Annealing Experiments

	Test Error %
LSTM+MLP(Reference)	33.28
Noisy LSTM+MLP( <b>NAN</b> )	31.12
Curriculum LSTM+MLP	14.83
Noisy LSTM+MLP( <b>NAN</b> ) Annealed Noise	<b>9.53</b>
Noisy LSTM+MLP( <b>NANIL</b> ) Annealed Noise	20.94

### NTM Experiments

