

- Activation functions are a crucial aspect of neural architectures, but are relatively underexplored in NAS research
- Previous work uses black-box optimization and thus comes with a high computational cost [1]

Motivation

- One-shot methods have offered a huge boost in efficiency in NAS research, which represent search spaces as a supergraph of architectures
- **Idea:** We apply Gradient based one-shot methods to search for novel activation functions and evaluate the results with respect to their performance and transferability

Technical Approach

Experimental Setup

Dataset: CIFAR-10

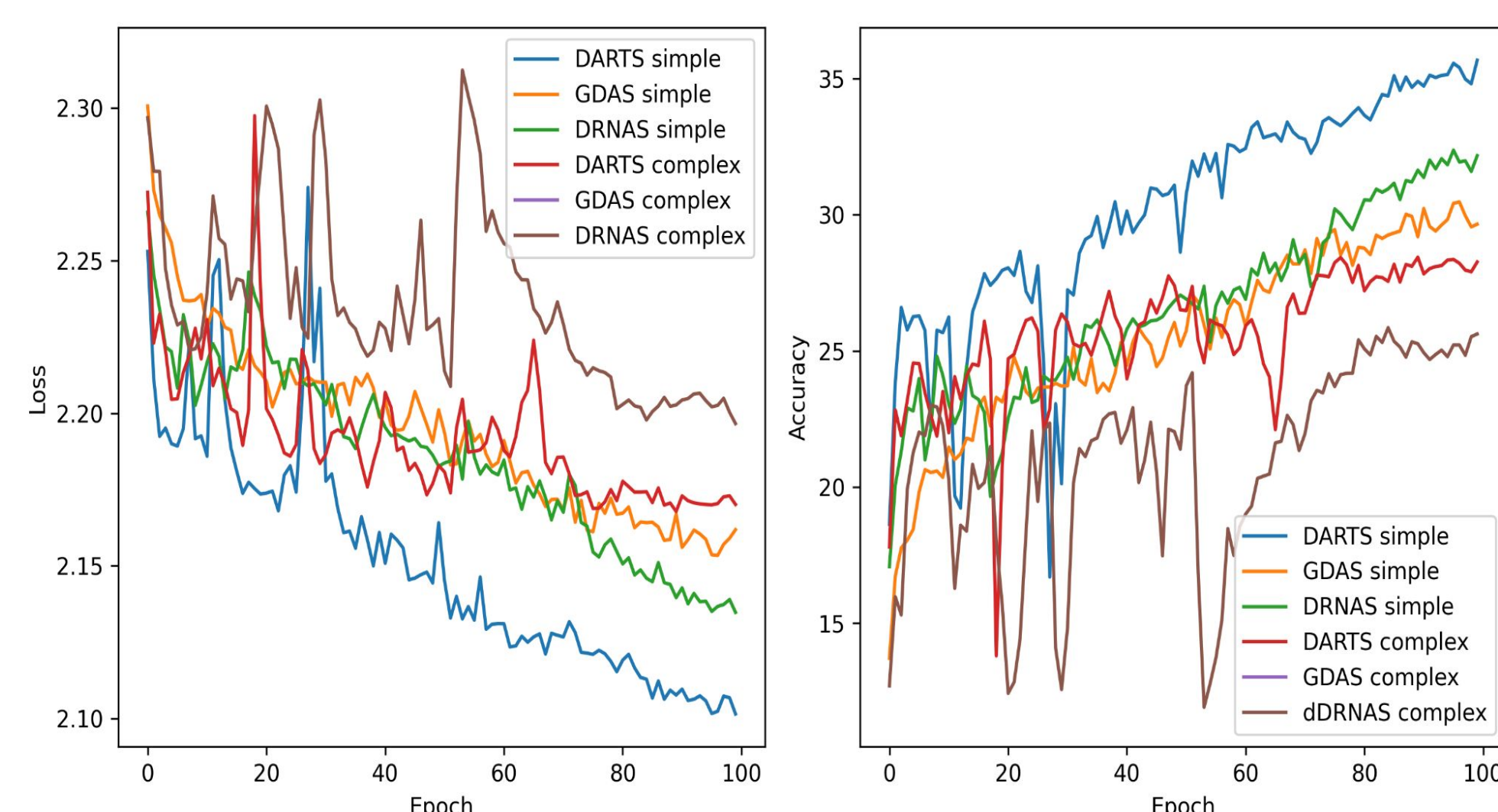
Optimizers: DARTS[2], GDAS[3], DRNAS[4]

Benchmarks: ReLU, Swish

Macro-architectures: ResNet8, ResNet20

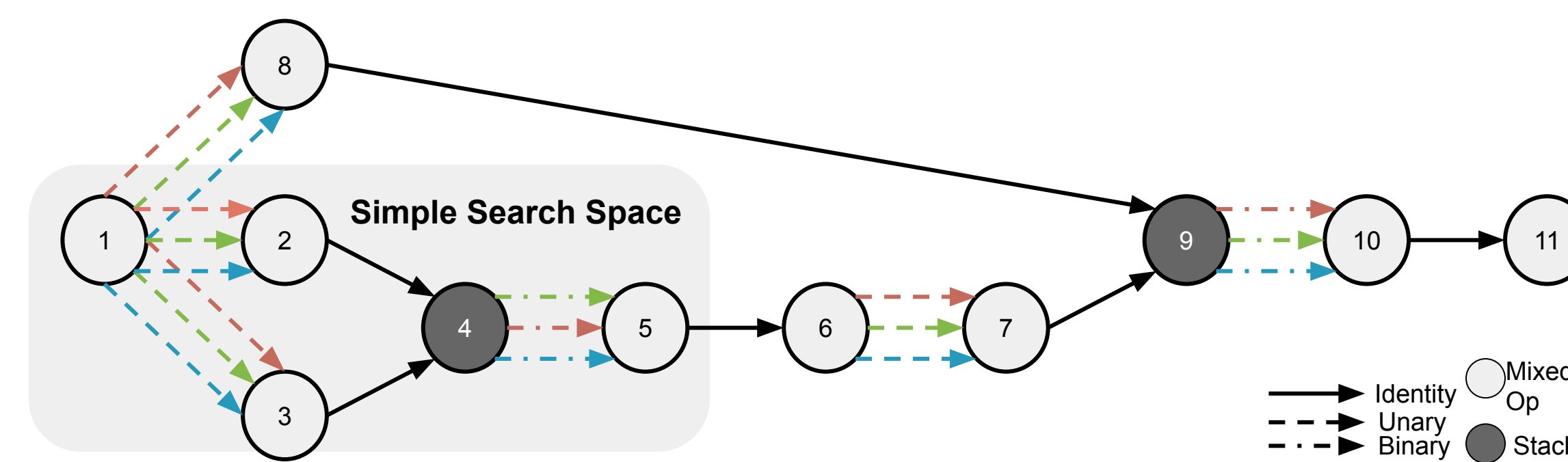
Unary Operators: x , $-x$, $|x|$, \sqrt{x} , βx , $x + \beta$, $\log(x)$, $\sin(x)$, $\cos(x)$, $\tanh(x)$, $\sinh^{-1}(x)$, $\tan^{-1}(x)$, $\text{sinc}(x)$, $\max(x, 0)$, $\min(x, 0)$, $\sigma(x)$, β

Binary Operators: $x_1 + x_2$, $x_1 \cdot x_2$, $x_1 - x_2$, $\max(x_1, x_2)$, $\min(x_1, x_2)$, $\sigma(x_1) \cdot x_2$, $\beta x_1 + (1 - \beta)x_2$ (β is a trainable parameter)



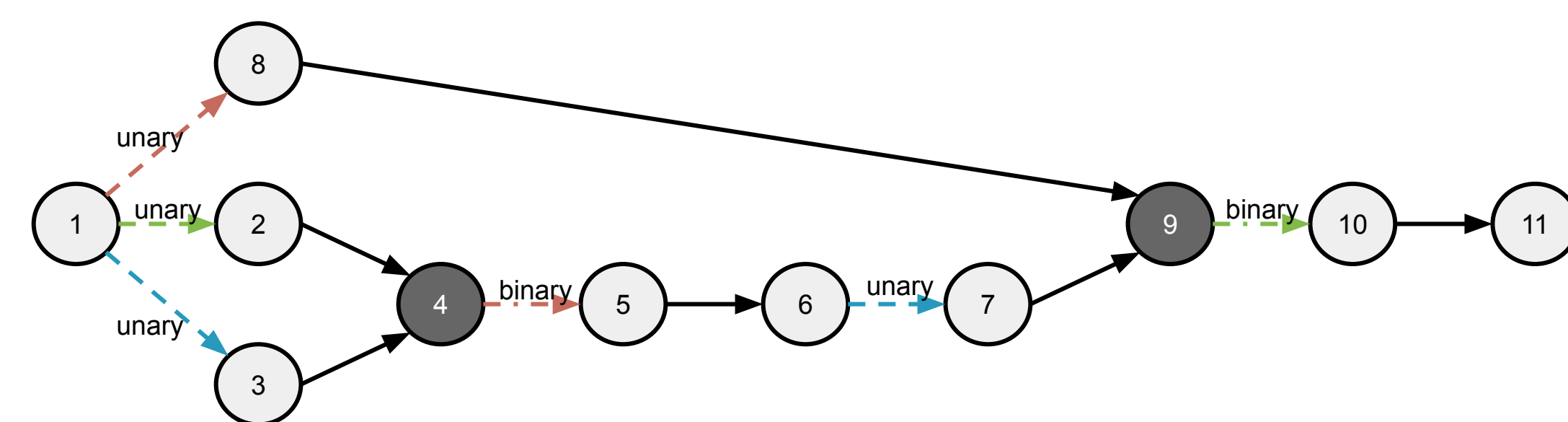
Search Space

- Using NASLib[5], we represent activation functions as cell-search-spaces of unary and binary operations as seen in [5]



- Then we insert these cells into ResNet architectures and optimize both the network and activation cell weights (alphas)

- Discretized network:

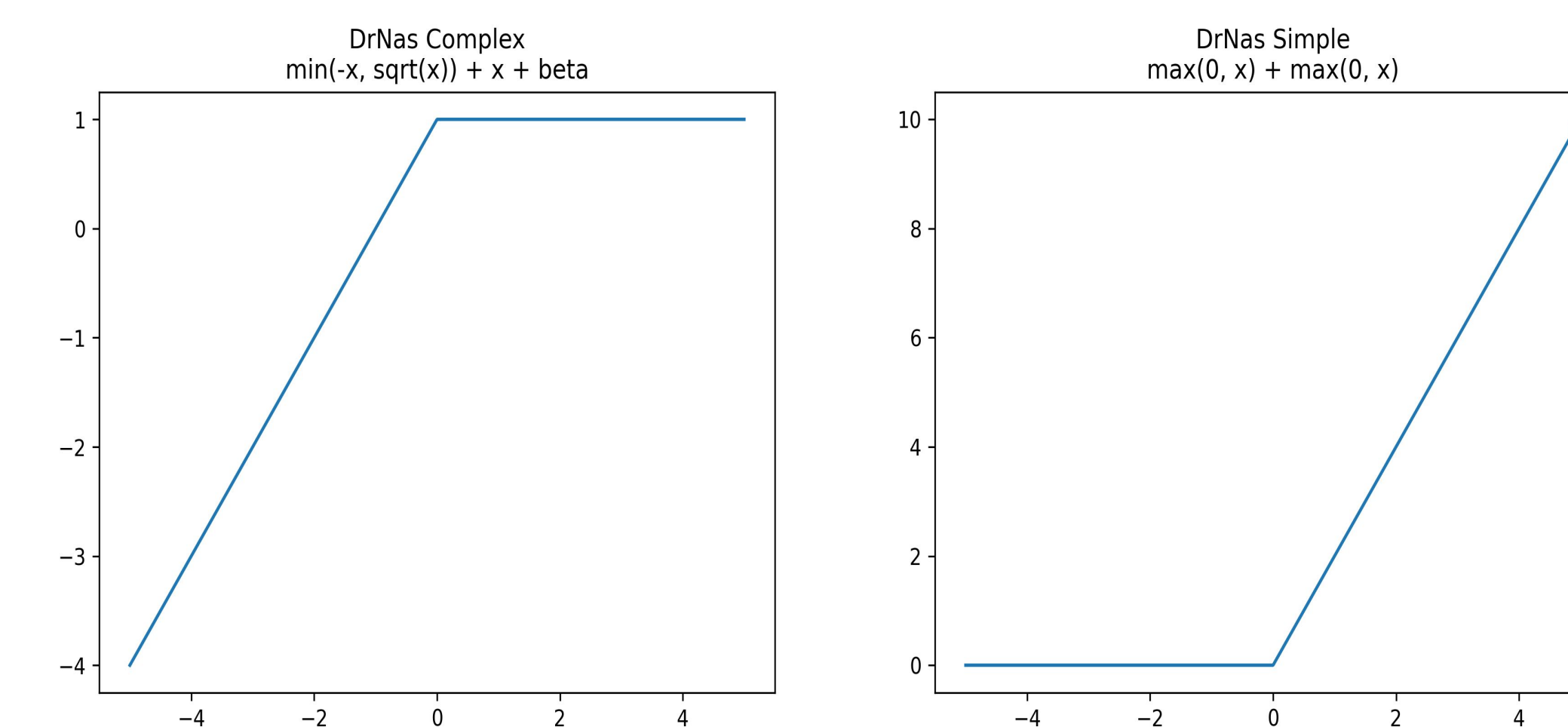
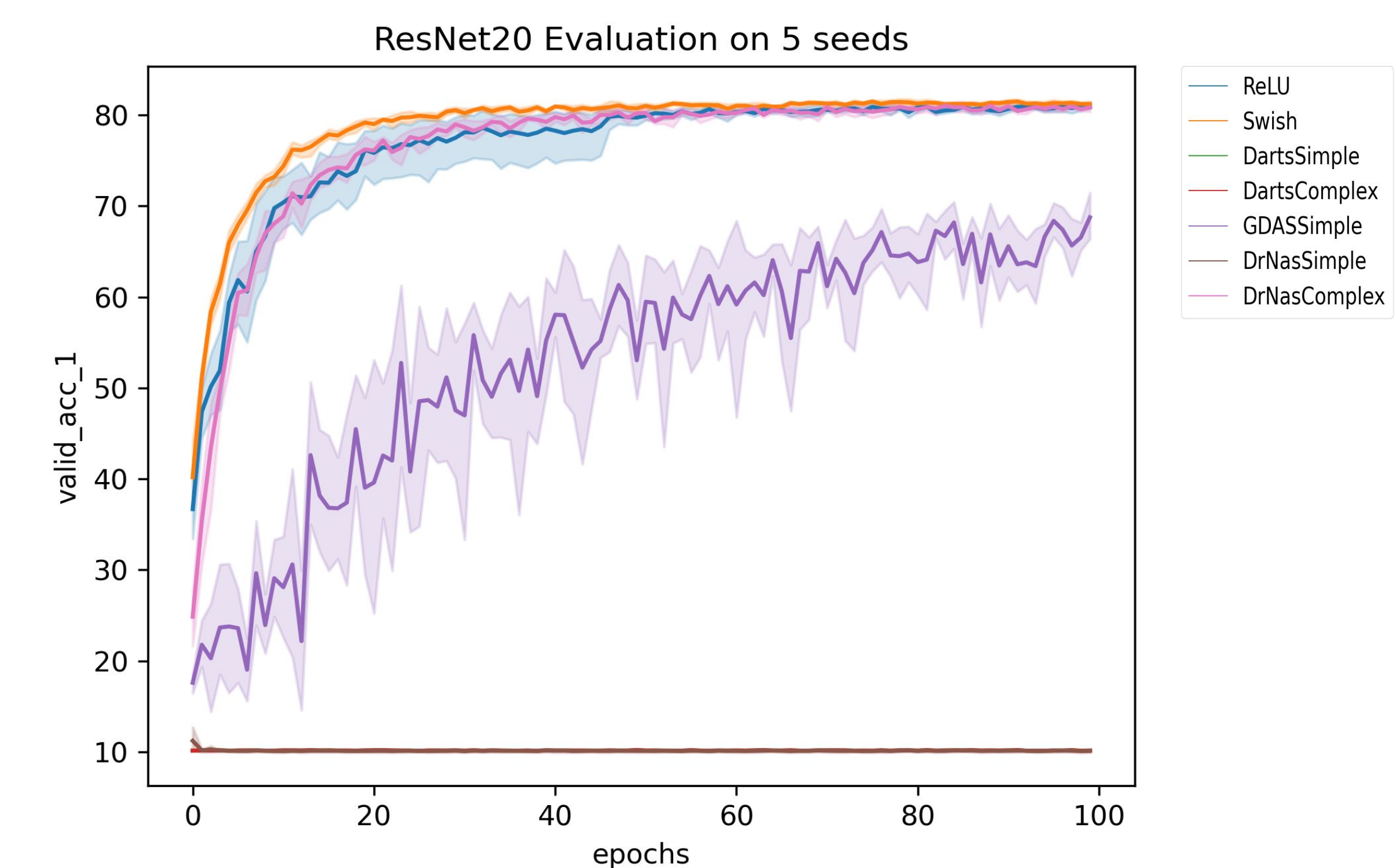
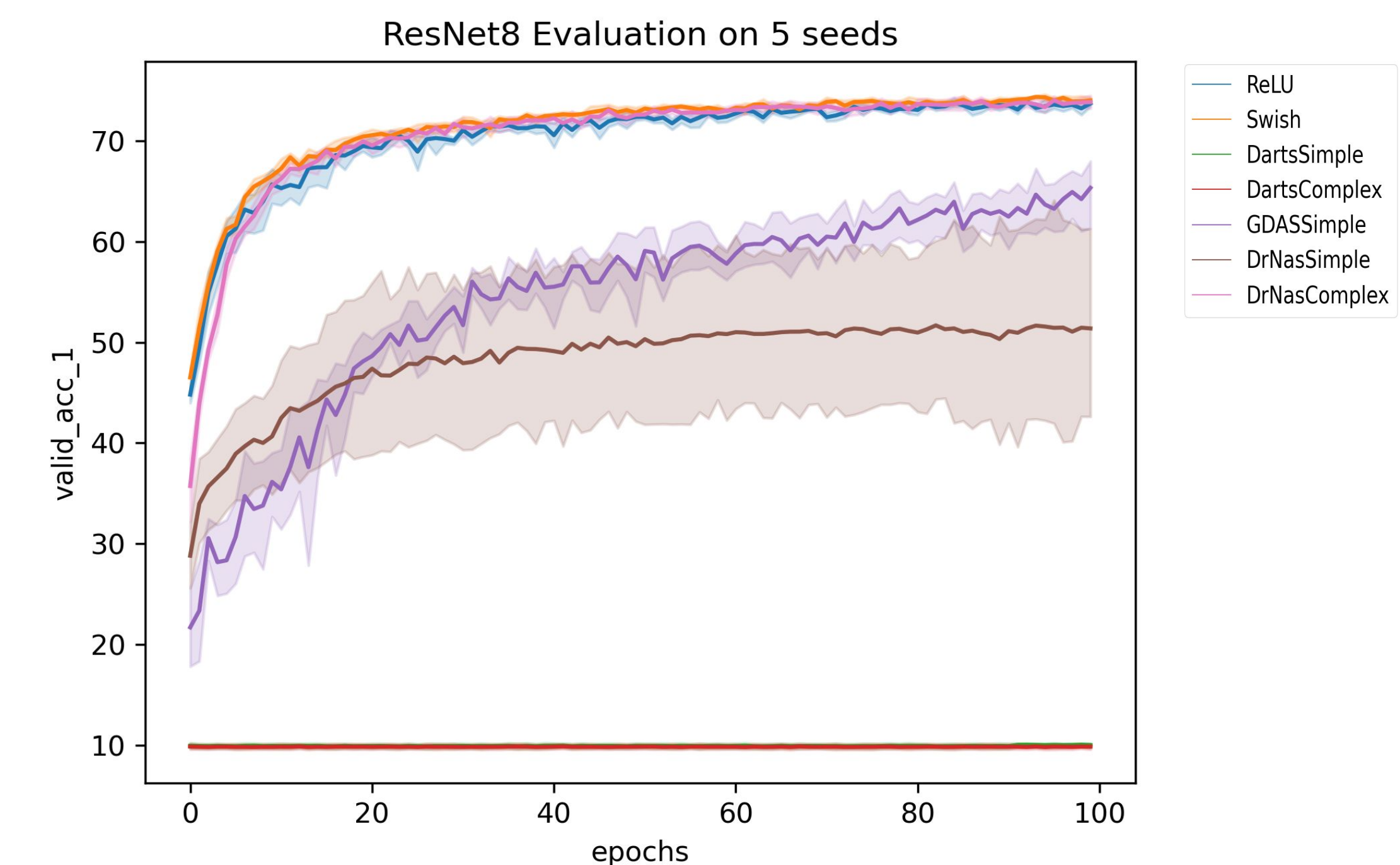


Summary/Conclusion

- The best performing activation is very similar to ReLU.
- There are limitations of DARTS based methods that affect how expressive the search space can be.
- Ensemble behavior of DARTS based methods affects the searched function. After discretization some functions fail to learn.
- Gradient based search will require additional modifications to the algorithm to account for unique behavior of activations and a carefully designed search space to produce novel activation function.

Experimental Results

- DrNas was able to find a function that has comparable performance to ReLU and Swish
- The functions found by DrNas are similar to ReLU ($\text{ReLU}(x) \cdot 2$ and $\text{Min}(x, 0) + \beta$).
- The training time for the DrNas activation is higher than ReLU and Swish without giving any additional performance increase
- Search loss and accuracy from Search search does not translate to evaluation.
- DARTS searched activations fail to learn in evaluation.
- Transferability across models seems to hold generally but the DrNasSimple activation fails to learn on the ResNet20
- Using operations such as $\exp(x)$ or $\text{Pow}(x)$ which have large values and gradients along with small operators adversely affects the Oneshot methods.



Activation Function		ResNet8 Test Accuracy		ResNet20 Test Accuracy	
		mean	std	mean	std
DartsComplex	$(x ^2 + \beta) + \beta x$	50.01	0.05	49.97	0.00
DartsSimple	$(\beta + x)^2$	50.01	0.05	49.97	0.00
DrNasComplex	$\min(-x, \sqrt{x}) + \beta + x$	97.76	0.11	98.42	0.15
DrNasSimple	$\max(0, x) + \max(0, x)$	79.74	8.07	49.98	0.02
GDASSimple	$\sigma(\sigma(x)) \cdot \sigma(x)$	88.58	5.54	94.96	1.13
ReLU	$\max(0, x)$	97.88	0.20	98.40	0.14
Swish	$x \cdot \sigma(x)$	97.76	0.14	98.48	0.11