

System Architecting Approach for Designing Deep Neural Networks

Cihan H Dagli

Missouri University of Science and Technology



System Architecting Approach for Designing Deep Neural Networks

Ram Deepak Gottapu
Cihan H Dagli

Need for this approach

- Manual effort
- Hyper-tuning of parameters
- Computation



Objective Function

- General objective function

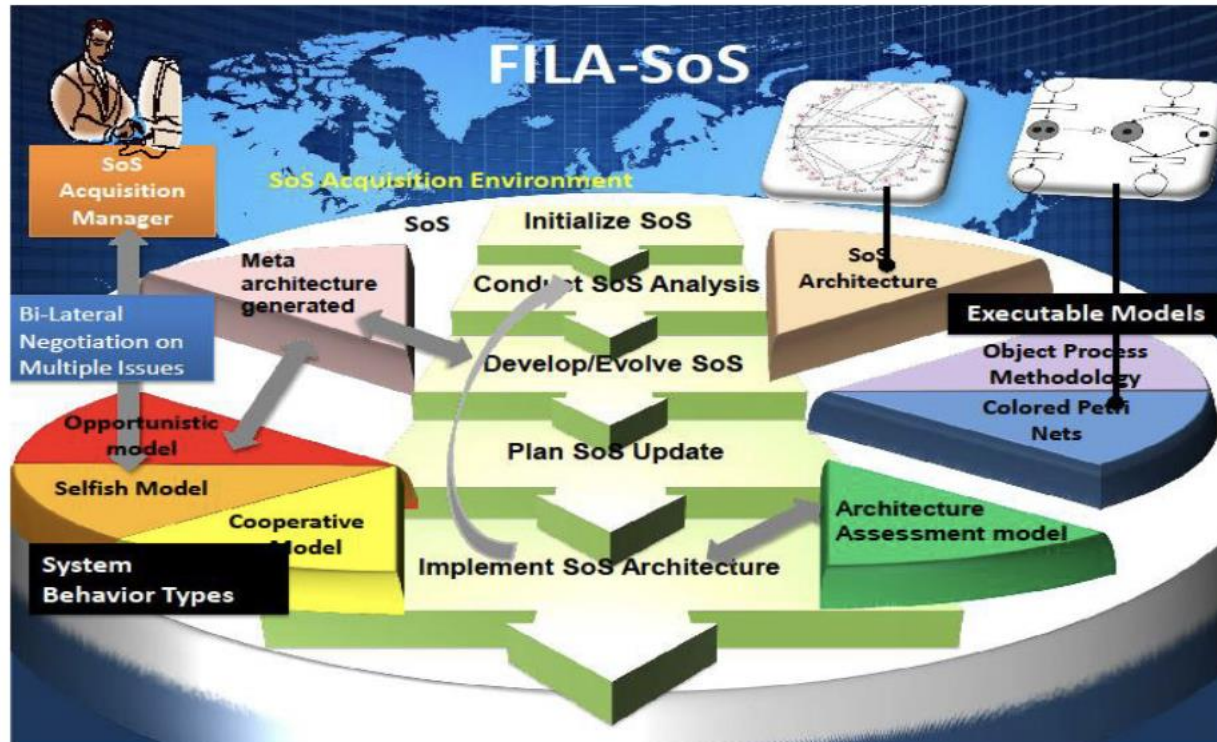
$$\underset{w}{\operatorname{argmin}} L_{\text{train}}(w, \alpha)$$

- Updated objective function

$$\underset{\alpha}{\min} L_{\text{val}}(w^*(\alpha), \alpha)$$

$$s.t. \quad w^*(\alpha) = \underset{w}{\operatorname{argmin}} L_{\text{train}}(w, \alpha)$$

Origin



Sos Approach

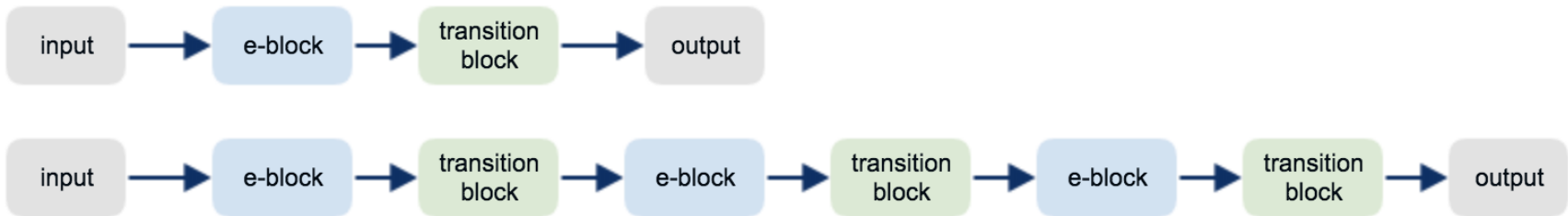
- Choose multiple objectives.
- Choose search space
- Can be used for both CNN and LSTM

Systems and Capabilities

System/Layer	Linear transformation	Non-linear transformation	Regularization	Parameter scaling
convolutional layer	X	-	-	-
pooling layer	X	-	-	-
activation layer	-	X	-	-
dropout	-	-	X	-
batch-normalization	-	-	X	X

Pre-determined design

- Reduce search space
- Stack blocks instead of layers



chromosome

- Select filter size
- Select stride
- Select compression
- Select inputs between h_0 and h_{t-1}
- Length of chromosome: $n_{total} d + (d(d + 1)/2)$

0	1	1	1	1	1	0	0	0
choice 1	choice 2				choice n		connections			

Search space

layer	choices
conv	1,3,5,7
Number of filters	12,24,36,48
compression	0.5, 0.65, 0.75, 1
stride	1,2

Evolution

- Selection
- Crossover
- Mutation
- To prevent training repeated architectures
 - Remove the possibility of generating identical chromosomes both in crossover and mutation by keeping a log of all the chromosomes generated from the beginning.
 - Add randomness to generate few random chromosomes in each generation to prevent the algorithm being struck in local minima.

Training conditions

- Dataset: CIFAR 10 (60000 labelled images with 10 classes)
- Optimization: stochastic gradient descent (SGD)
- learning rate: 0.1
- batch size: 64
- Total duration: 1 month
- No. of architectures searched: approx 1000

Limitations

- Size of each e-block is pre-determined
- Different size e-block requires re-training
- Computationally expensive

results

Method	Depth	<u>Params (M)</u>	CIFAR-10 error percentage
Network in Network [6]	-	-	8.81
All CNN [7]	-	-	7.25
Deeply Supervised Net [5]	-	-	7.97
Fractal Net [4]	21	38.6	5.22
<u>ResNet</u> [10]	110	1.7	6.61
Wide <u>ResNet</u> [8]	16	11	4.81
<u>ResNet</u> (pre-activation) [2]	164	1.7	5.46
<u>DenseNet</u> [3]	100	0.8	4.51
e-block@1	100	1.2	5.23
e-block@2	100	2.8	4.61