

**Benchmarks:** All benchmark tests are run with the same parameters over 10 randomly selected seeds. *We decide to train all methods until they converge, and this is how we choose how many iterations to run each of them for.*

## FMNIST:

**PopDescent:** All parameters are the same for the without/with regularization tests, except for changing the model by adding one l2 kernel regularizer to the second to last fully connected layer in the model with regularization.

### Training Parameters:

<i>Population Size</i>	5
<i>Replaced Individuals (in m-elitist)</i>	2
<i>Iterations</i>	50
<i>Batches</i>	128
<i>Batch Size</i>	64
<i>lr</i>	0.001
<i>Regularization Rate (in model with regularization)</i>	0.001
<i>Optimizer</i>	Adam

### Randomization Parameters:

<i>lr_constant</i>	$10^{**}(\text{random.normal}(\mu=-4, \sigma=2))$
<i>regularization_constant</i>	$10^{**}(\text{random.normal}(\mu=0, \sigma=2))$
<i>randomization_amount</i> (amount to randomize model, changes based on model's performance)	$1 - (2 / (2 + \text{model loss}))$
<i>Gaussian Noise for Model Weights</i> (sum of current weights and noise)	$\text{noise} = \text{random.normal}(\mu=0, \sigma=0.01) * \text{randomization\_amount}$
<i>Gaussian Noise for lr_constant</i> (product of current lr_constant and noise)	$2^{**}(\text{np.random.normal}(\mu=0, \sigma=\text{randomization\_amount}*15))$
<i>Gaussian Noise for regularization_rate</i> (product of current regularization_constant and noise)	$2^{**}(\text{np.random.normal}(\mu=0, \sigma=\text{randomization\_amount}*15))$

**Basic Grid Search:** All parameters are the same for the without/with regularization tests. Both without/with regularization tests use the model with l2 kernel regularization, but when training without regularization, the *Regularization Rates to Sample* is simply [0], meaning no regularization.

**Training Parameters:**

<i>Learning Rates to Sample</i>	[0.01, 0.001, 0.0001, 0.00001, 0.000001]
<i>Regularization Rates to Sample</i>	[0.01, 0.001, 0.0001, 0.00001, 0.000001]
<i>Iterations</i>	100
<i>Batches</i>	128
<i>Batch Size</i>	64
<i>Optimizer</i>	Adam

**KT RandomSearch:** All parameters are the same for the without/with regularization tests, except for changing the model by adding one l2 kernel regularizer to the second to last fully connected layer in the model with regularization. We used Keras Tuner's RandomSearch implementation.

**Training Parameters:**

<i>Range of Learning Rates to Sample (continuous)</i>	[0.0001, 0.01]
<i>Range of Regularization Rates to Sample</i>	[0.00001, 0.1]
<i>max_trials</i>	25
<i>executions_per_trial</i>	2
<i>train_epochs (with early stopping, patience=2; train over the whole train dataset)</i>	20
<i>Batch Size</i>	64
<i>Optimizer</i>	Adam

**ESGD:** We used *tqch*'s open-source implementation of ESGD on <https://github.com/tqch/esgd-ws>.

**Training Parameters:**

<i>n_population</i>	5
<i>sgds_per_gen</i>	1

<i>evos_per_gen</i>	1
<i>reproductive_factor</i>	4
<i>m_elite</i>	3
<i>mixing_number</i>	3
<i>optimizer_class</i>	SGD (lr=0.001)
<i>n_generations</i>	10
<i>batch_size</i>	64

## CIFAR-10:

**PopDescent:** All parameters are the same for the without/with regularization tests, except for changing the model by adding one l2 kernel regularizer to the second to last fully connected layer in the model with regularization.

### Training Parameters:

<i>Population Size</i>	5
<i>Replaced Individuals (in m-elitist)</i>	2
<i>Iterations</i>	20
<i>Batches (trained over 2 epochs per iteration here)</i>	128
<i>Batch Size</i>	64
<i>lr</i>	0.001
<i>Regularization Rate (in model with regularization)</i>	0.001
<i>Optimizer</i>	Adam

### Randomization Parameters:

<i>lr_constant</i>	$10^{**}(\text{random.normal}(\mu=-4, \sigma=2))$
<i>regularization_constant</i>	$10^{**}(\text{random.normal}(\mu=0, \sigma=2))$
<i>randomization_amount</i> (amount to randomize model, changes based on model's performance)	$1 - (2 / (2 + \text{model loss}))$
<i>Gaussian Noise for Model Weights (sum of current</i>	$\text{noise} = \text{random.normal}(\mu=0,$

<i>weights and noise</i> )	$\text{sigma}=0.01) * \text{randomization\_amount}$
<i>Gaussian Noise for lr_constant (product of current lr_constant and noise)</i>	$2^{**}(\text{np.random.normal}(\mu=0, \text{sigma}=\text{randomization\_amount} * 15))$
<i>Gaussian Noise for regularization_rate (product of current regularization_constant and noise)</i>	$2^{**}(\text{np.random.normal}(\mu=0, \text{sigma}=\text{randomization\_amount} * 15))$

**Basic Grid Search:** All parameters are the same for the without/with regularization tests. Both without/with regularization tests use the model with l2 kernel regularization, but when training without regularization, the *Regularization Rates to Sample* is simply [0], meaning no regularization.

**Training Parameters:**

<i>Learning Rates to Sample</i>	[0.01, 0.001, 0.0001, 0.00001, 0.000001]
<i>Regularization Rates to Sample</i>	[0.01, 0.001, 0.0001, 0.00001, 0.000001]
<i>Iterations</i>	30
<i>Batches</i>	128
<i>Batch Size</i>	64
<i>Optimizer</i>	Adam

**KT RandomSearch:** All parameters are the same for the without/with regularization tests, except for changing the model by adding one l2 kernel regularizer to the second to last fully connected layer in the model with regularization. We used Keras Tuner's RandomSearch implementation.

**Training Parameters:**

<i>Range of Learning Rates to Sample (continuous)</i>	[0.0001, 0.01]
<i>Range of Regularization Rates to Sample</i>	[0.00001, 0.1]
<i>max_trials</i>	25
<i>executions_per_trial</i>	2
<i>train_epochs (with early stopping, patience=2; train over the whole train dataset)</i>	20
<i>Batch Size</i>	64
<i>Optimizer</i>	Adam

**ESGD:** We used *tqch*'s open-source implementation of ESGD on <https://github.com/tqch/esgd-ws>.

**Training Parameters:**

<i>n_population</i>	5
<i>sgds_per_gen</i>	1
<i>evos_per_gen</i>	1
<i>reproductive_factor</i>	4
<i>m_elite</i>	3
<i>mixing_number</i>	3
<i>optimizer_class</i>	SGD (lr=0.001)
<i>n_generations</i>	3
<i>batch_size</i>	8

***Convergence Test:*** All convergence tests are run with the same parameters over 6 randomly selected seeds. All tests are run on the same model without regularization on the FMNIST dataset.

**PopDescent:**

**Training Parameters**

<i>Population Size</i>	5
<i>Replaced Individuals (in m-elitist)</i>	2
<i>Iterations</i>	115
<i>Batches</i>	128
<i>Batch Size</i>	64
<i>lr</i>	0.001
<i>Optimizer</i>	Adam

**Randomization Parameters:**

<i>lr_constant</i>	10** (random.normal(mu=-4, sigma=2))
--------------------	--------------------------------------

<i>randomization_amount</i> (amount to randomize model, changes based on model's performance)	$1 - (2 / (2 + \text{model loss}))$
<i>Gaussian Noise for Model Weights</i> (sum of current weights and noise)	<code>noise = random.normal(mu=0, sigma=0.01)*randomization_amount</code>
<i>Gaussian Noise for lr_constant</i> (product of current lr_constant and noise)	<code>2**(np.random.normal(mu=0, sigma=randomization_amount*15))</code>

## Basic Grid Search:

### Training Parameters

<i>Learning Rates to Sample</i>	[0.01, 0.001, 0.0001, 0.00001, 0.000001]
<i>Regularization Rates to Sample</i>	[0]
<i>Iterations</i>	100
<i>Batches</i>	128
<i>Batch Size</i>	64
<i>Optimizer</i>	Adam

## KT RandomSearch:

### Training Parameters

<i>Range of Learning Rates to Sample (continuous)</i>	[0.0001, 0.01]
<i>Range of Regularization Rates to Sample</i>	[0.00001, 0.1]
<i>max_trials</i>	25
<i>executions_per_trial</i>	2
<i>train_epochs</i> (with early stopping, patience=2; train over the whole train dataset)	20
<i>Batch Size</i>	64
<i>Optimizer</i>	Adam

## ESGD:

**Training Parameters**

<i>n_population</i>	5
<i>sgds_per_gen</i>	1
<i>evos_per_gen</i>	1
<i>reproductive_factor</i>	4
<i>m_elite</i>	3
<i>mixing_number</i>	3
<i>optimizer_class</i>	SGD (lr=0.001)
<i>n_generations</i>	15
<i>batch_size</i>	64

***Ablation Study:*** To emphasize the differences of PopDescent’s features, we train only on the first 10k images in the FMNIST dataset, and add l2 kernel regularization to every layer in the benchmark FMNIST model, initialized with a default value of 0.001 in each layer. All training parameters are the same in each of the four tests listed in the paper. We only change whether PopDescent’s randomization scheme is on or off, with CV selection turned on, and using the model with regularization for the first two tests listed (“Ablation Study Over Randomization”). We only change the selection process when comparing without CV selection vs with CV selection, both on a model without regularization for the second two tests listed (“Ablation Study Over Cross-Validation Fitness”).

**Training Parameters:**

<i>Population Size</i>	10
<i>Replaced Individuals (in m-elitist)</i>	5
<i>Iterations</i>	35
<i>Batches</i>	128
<i>Batch Size</i>	64
<i>lr</i>	0.001
<i>Regularization Rate</i>	0.001

<i>Optimizer</i>	Adam
------------------	------

#### Randomization Parameters:

<i>lr_constant</i>	$10^{**}(\text{random.normal}(\mu=-4, \sigma=2))$
<i>regularization_constant</i>	$10^{**}(\text{random.normal}(\mu=0, \sigma=2))$
<i>randomization_amount</i> (amount to randomize model, changes based on model's performance)	$1 - (2 / (2 + \text{model loss}))$
<i>Gaussian Noise for Model Weights</i> (sum of current weights and noise)	$\text{noise} = \text{random.normal}(\mu=0, \sigma=0.01) * \text{randomization\_amount}$
<i>Gaussian Noise for lr_constant</i> (product of current lr_constant and noise)	$2^{**}(\text{np.random.normal}(\mu=0, \sigma=\text{randomization\_amount} * 15))$
<i>Gaussian Noise for regularization_rate</i> (product of current regularization_constant and noise)	$2^{**}(\text{np.random.normal}(\mu=0, \sigma=\text{randomization\_amount} * 15))$

***Hyperparameter Sensitivity Study:*** We conduct this study on the CIFAR-10 dataset, on the same model used in the CIFAR-10 benchmark tests without regularization. When testing for the effects of changing the number of iterations vs different learning rates, we keep everything constant, except for changing the hyperparameter we are testing by sampling each possibility in the arrays listed for that category.

#### Training Parameters:

<i>Population Size</i>	10
<i>Replaced Individuals (in m-elitist)</i>	5
<b><i>Iterations</i></b>	<b>[10, 30, 50]</b>
<i>Batches</i>	128
<i>Batch Size</i>	64
<b><i>lr</i></b>	<b>[0.001, 0.01, 0.05]</b>
<i>Regularization Rate</i>	0.001
<i>Optimizer</i>	Adam



**Randomization Parameters:**

<i>lr_constant</i>	$10^{**}(\text{random.normal}(\mu=-4, \sigma=2))$
<i>regularization_constant</i>	$10^{**}(\text{random.normal}(\mu=0, \sigma=2))$
<i>randomization_amount</i> (amount to randomize model, changes based on model's performance)	$1 - (2 / (2 + \text{model loss}))$
<i>Gaussian Noise for Model Weights</i> (sum of current weights and noise)	$\text{noise} = \text{random.normal}(\mu=0, \sigma=0.01) * \text{randomization\_amount}$
<i>Gaussian Noise for lr_constant</i> (product of current lr_constant and noise)	$2^{**}(\text{np.random.normal}(\mu=0, \sigma=\text{randomization\_amount} * 15))$
<i>Gaussian Noise for regularization_rate</i> (product of current regularization_constant and noise)	$2^{**}(\text{np.random.normal}(\mu=0, \sigma=\text{randomization\_amount} * 15))$

**ESGD:****Training Parameters:**

<i>n_population</i>	5
<i>sgds_per_gen</i>	1
<i>evos_per_gen</i>	1
<i>reproductive_factor</i>	4
<i>m_elite</i>	3
<i>mixing_number</i>	3
<i>optimizer_class</i>	SGD
<i>learning_rate</i>	[0.001, 0.01, 0.05]
<i>n_generations</i>	[1, 3, 5]
<i>batch_size</i>	8