



RUB

RUHR-UNIVERSITÄT BOCHUM

# UNCERTAINTY ESTIMATION IN FEDERATED DEEP LEARNING

Florian Linsner

# Table of Contents

- Mathematical Basics
- Uncertainty Theory
- Methods for Uncertainty
- Challenges in Centralized Deep Learning
- Federated Deep Learning
- Leveraging Uncertainty in Federated Deep Learning
- Empirical Evaluation
- Conclusion
- Discussion

Chapter

# Mathematical Basics

# Entropy and Variance

- Shannon Entropy over output classes

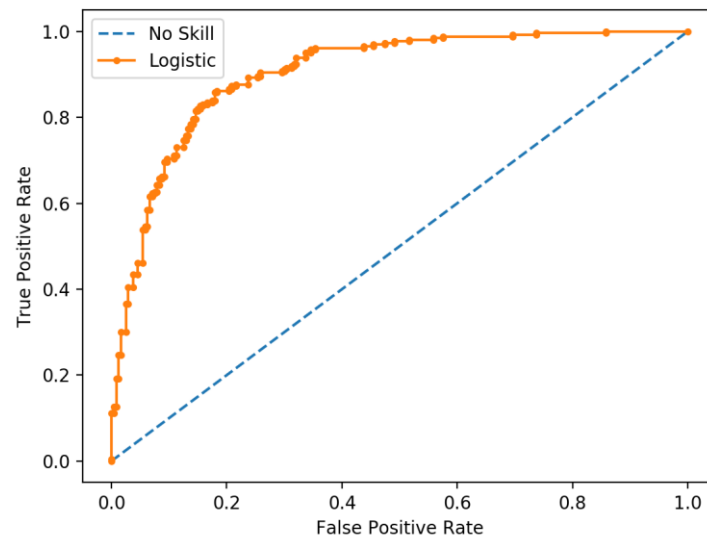
$$H(p(Y|x^*, \theta)) = - \sum_{k=1}^K p(y_k|x^*, \theta) \cdot \log p(y_k|x^*, \theta)$$

- Variance over multiple predictions for a fixed input

$$\sigma^2 = \frac{1}{S} \sum_{s=1}^S p(y|x^*, \theta_s)^2 - p(y|x^*)^2$$

# ROC and AUROC

- Receiver Operating Characteristic Curve
  - Binary classifier
  - 2 groups of data
  - True Positive vs False Positive
  - Multiple Thresholds
- Area Under ROC
  - Score between 0.5 and 1
  - Want to achieve a score of 1



Chapter

# Uncertainty Theorie

# Importance of Uncertainty

- We use models to make predictions
  - Can we trust these predictions?
- There is no perfect model
  - No-Free-Lunch-Theorem
  - Models make mistakes!
- Many Decisions in critical environments
  - Autonomous Driving
  - Nuclear Powerplants
  - Industrial Settings

# Malicious Intentions

- An Attacker can Influence the behaviour of a model
- Data Poisoning
- Aversarial Examples
- Many other

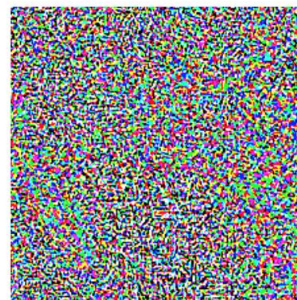


$x$

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

=



$x +$

$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence



# Bayesian Uncertainty - Origins

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\sum_{y \in Y} p(x|y)p(y)} = \frac{p(x|y)p(y)}{Z}$$


# Bayesian Uncertainty - Origins

Prior

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\sum_{y \in Y} p(x|y)p(y)} = \frac{p(x|y)p(y)}{Z}$$

# Bayesian Uncertainty - Origins

More Evidence      Prior


$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\sum_{y \in Y} p(x|y)p(y)} = \frac{p(x|y)p(y)}{Z}$$

# Bayesian Uncertainty - Origins

More Evidence      Prior

Posterior

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\sum_{y \in Y} p(x|y)p(y)} = \frac{p(x|y)p(y)}{Z}$$

# Bayesian Uncertainty - Origins

More Evidence

Prior

Posterior

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\sum_{y \in Y} p(x|y)p(y)} = \frac{p(x|y)p(y)}{Z}$$

Denominator Usually  
intractable

# Bayesian Uncertainty - Usage

- Stating priors over model parameters  $\beta$
- Using training data  $(X, y)$

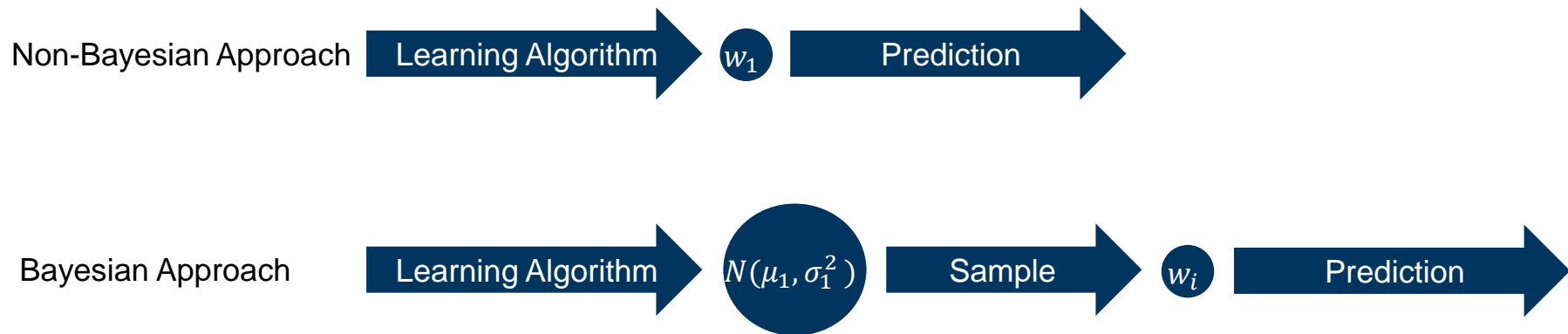
# Bayesian Uncertainty - Usage

- Stating priors over model parameters  $\beta$
- Using training data  $(X, y)$

$$p(\beta|y_t, X_t) = \frac{p(y_t, X_t|\beta)p(\beta)}{p(y_t, X_t)} = \frac{p(y_t, X_t|\beta)p(\beta)}{\int p(y_t, X_t|\beta)p(\beta)d\beta}$$

- Solving analytically is often impossible/impractical
- Solution: Numeric Approximations

# Bayesian Uncertainty - Usage



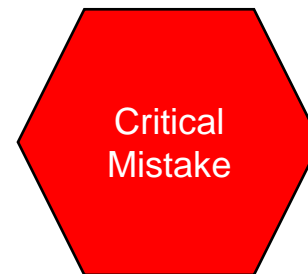
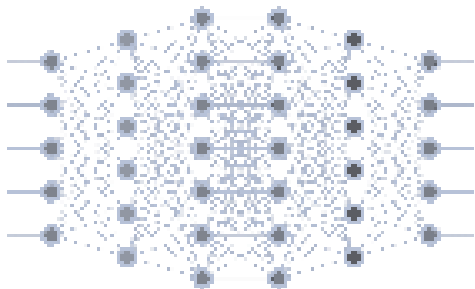


Chapter

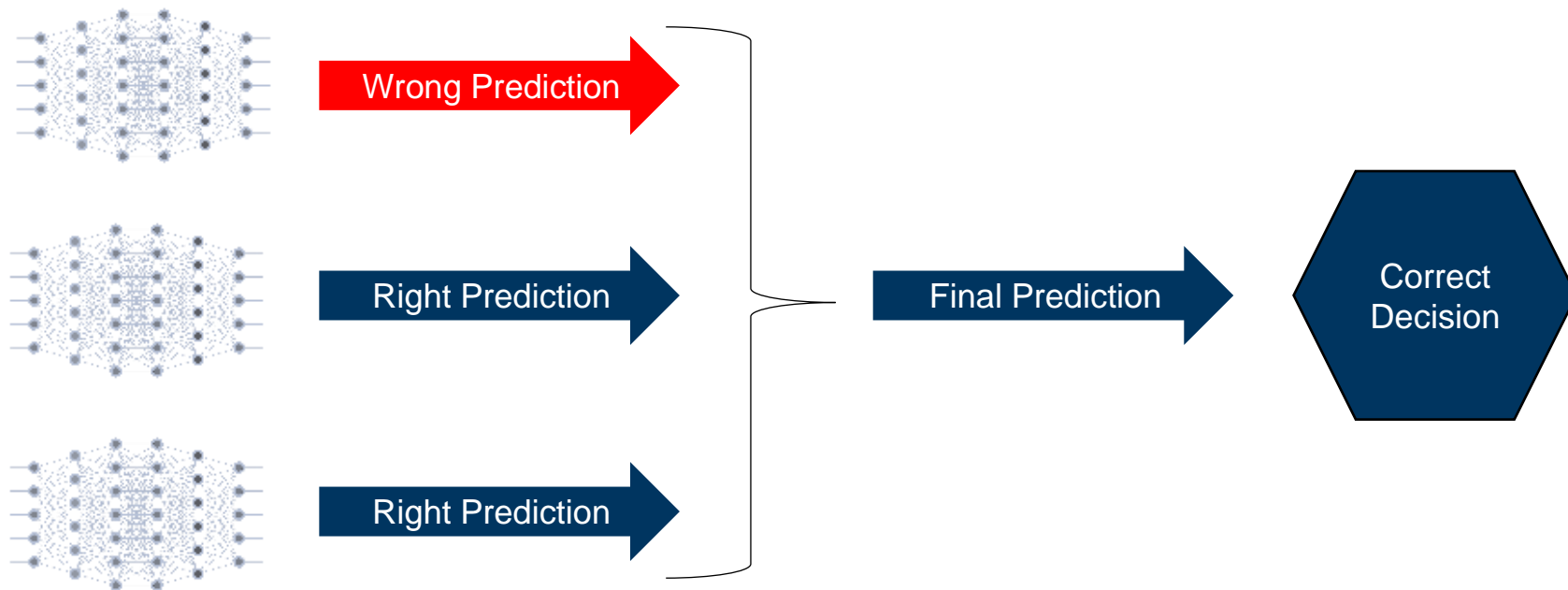
**Methods for Uncertainty**

# Ensembles

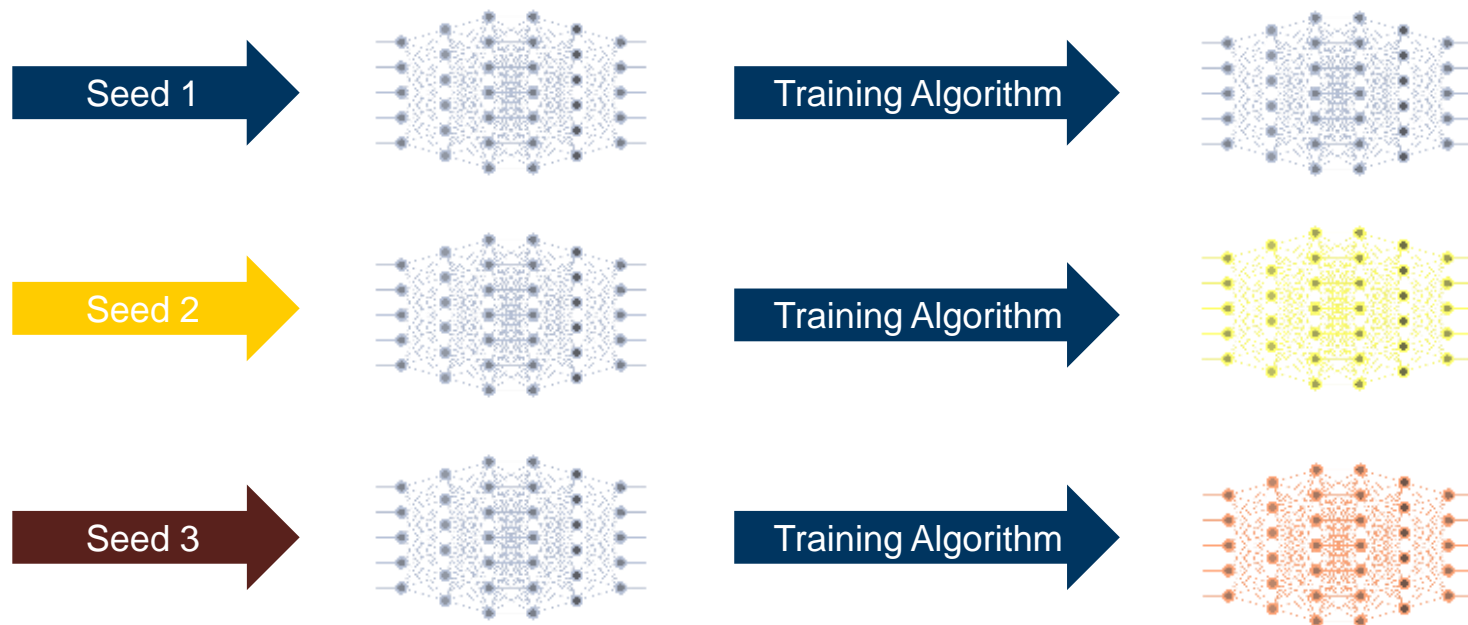
- Originally proposed as non-Bayesian method
- Can be viewed as Bayesian Approximation[1]



# Ensembles

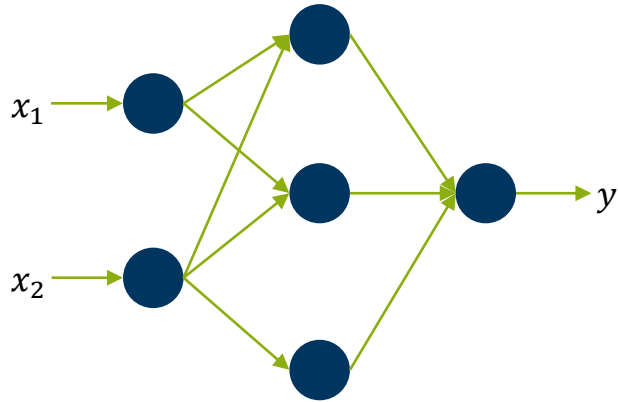


# Ensemble Training



# Dropout

- Method for Regularization during Training Time
- Prevents Overfitting



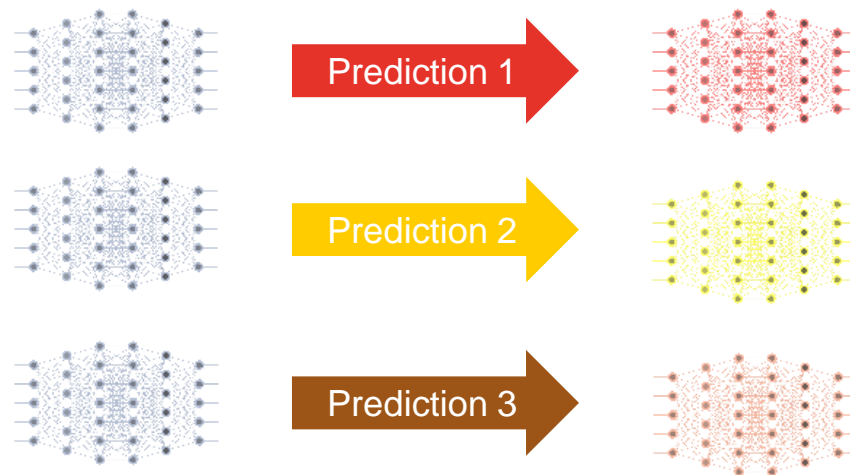
# Dropout

- Method for Regularization during Training Time
- Prevents Overfitting



# MC-Dropout

- MC-Dropout as Bayesian Approximation[2]
- Dropout During Prediction Time
- Multiple Predictions for same Input



# SWAG

- Stochastic Weight Averaging Gaussian[3]
- Needs Pre-Trained Solution near a local minima
- Compute a running average

$$\theta_{\text{SWA}} = \frac{1}{T} \sum_{i=1}^T \theta_i$$

- Compute squared running average

$$\overline{\theta^2} = \frac{1}{T} \sum_{i=1}^T \theta_i^2$$



# SWAG

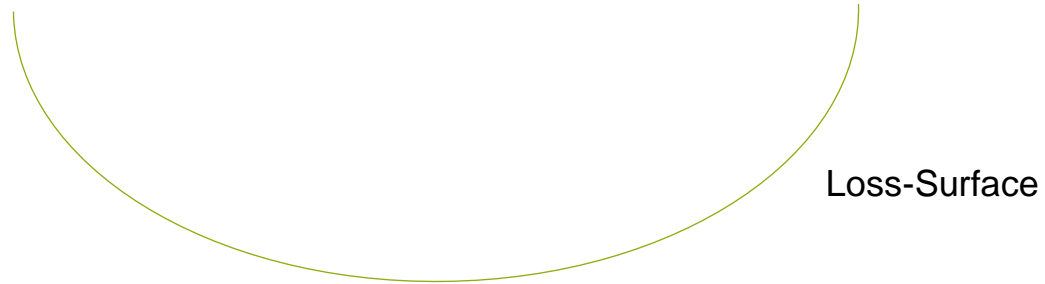
- Compute Covariance

$$\Sigma_{\text{diag}} = \text{diag}(\overline{\theta^2} - \theta_{\text{SWA}}^2)$$

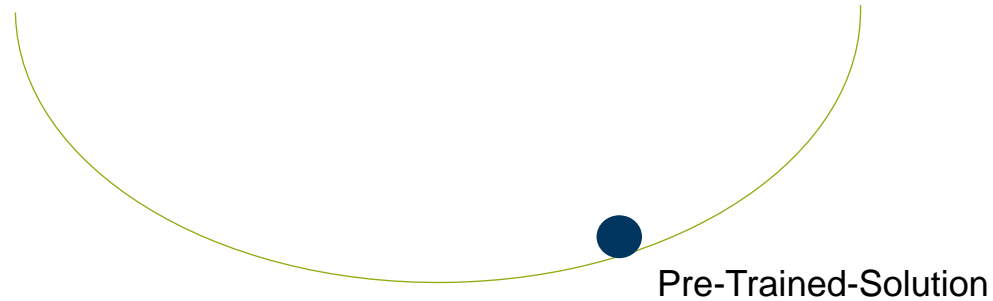
- Sample from Distribution

$$\mathcal{N}(\theta_{\text{SWA}}, \Sigma_{\text{Diag}})$$

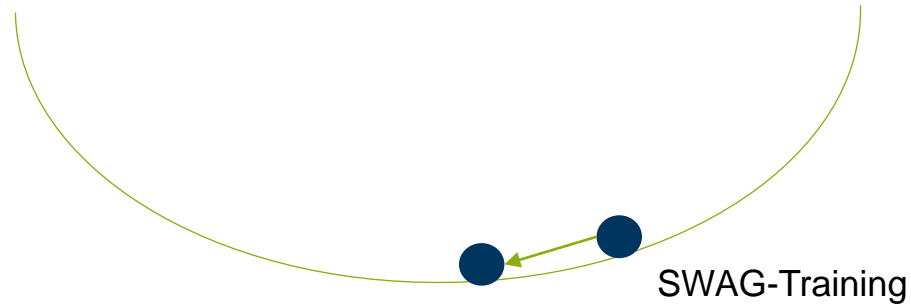
# SWAG in Pictures



# SWAG in Pictures

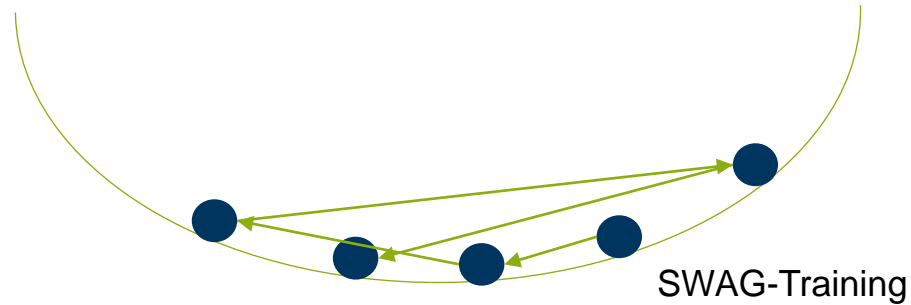


# SWAG in Pictures



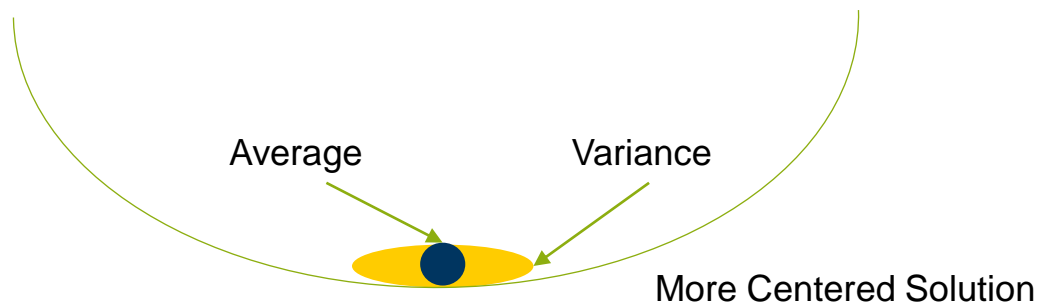
# SWAG in Pictures

- Goal: More Centered Solution



# SWAG in Pictures

- Goal: More Centered Solution



Chapter

**Challenges in Centralized  
Deep Learning**

# Data Generating Devices



- $50.1 \cdot 10^9$  connected devices by 2020[4]
  - (autonomous) driving and vehicles
  - Smart sensors
  - Mobile phones
  - Airplanes
  - Entertainment Devices
  - Industrial Machines
  - ...
- IoT sector is growing exponentially



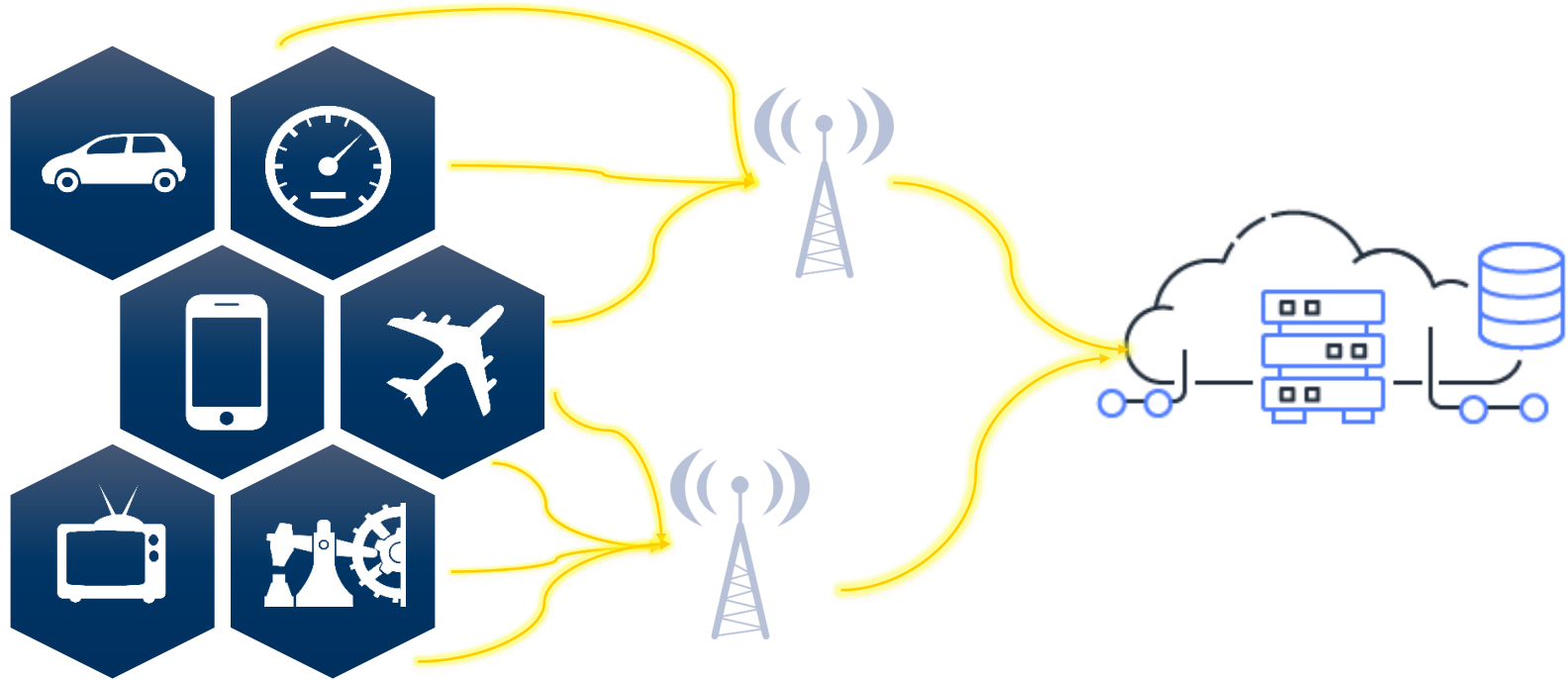
# Centralized Machine Learning



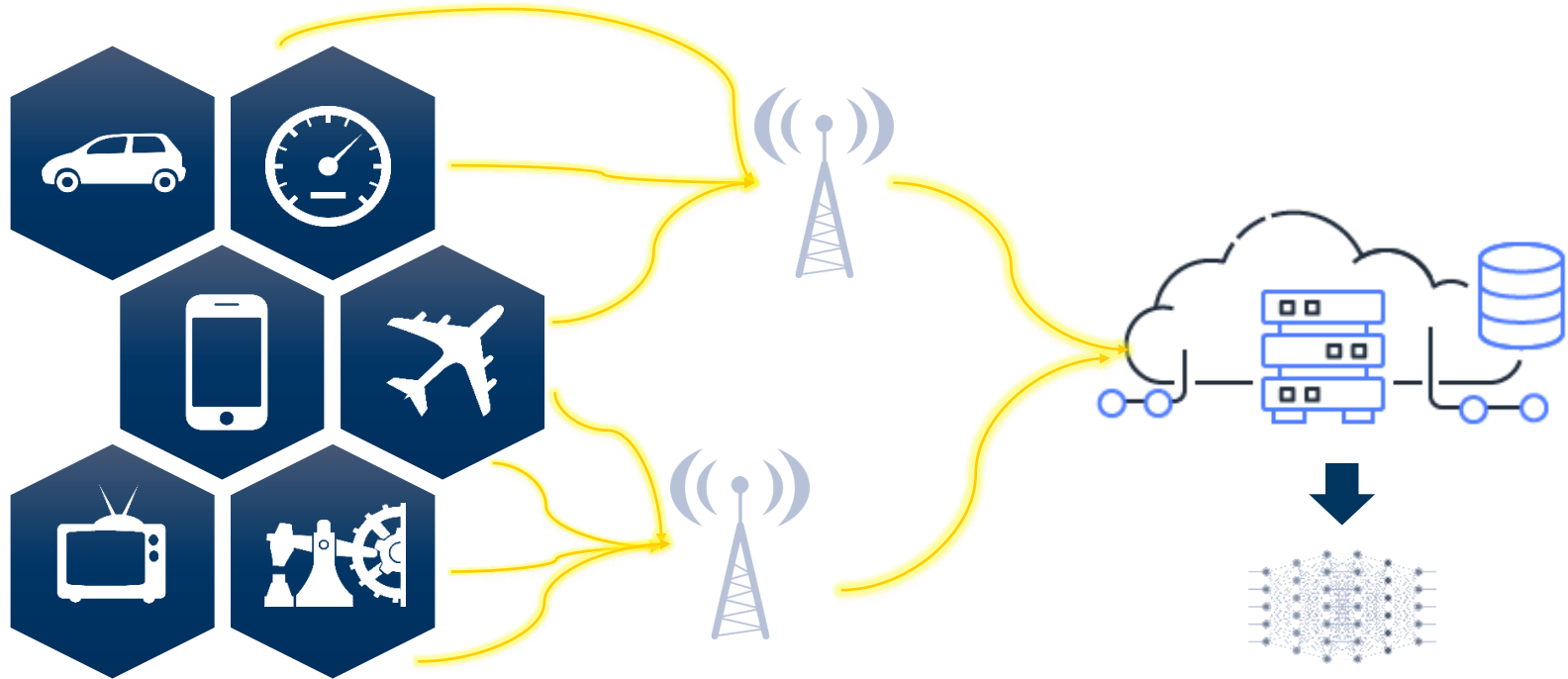
# Centralized Machine Learning



# Centralized Machine Learning



# Centralized Machine Learning



# Limits of Cloud Processing



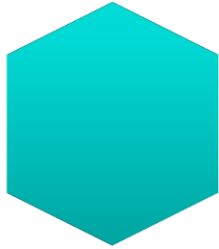
- Autonomous vehicle generates 1GB per second[5]
  - 4G mobile: 37.5 MB per second
  - VW sold  $10.8 \cdot 10^6$  cars in 2018
  - Produce 10.3 PB data per second
  - All CERN experiments process 25 GB per second[6]

# Limits of Cloud Processing



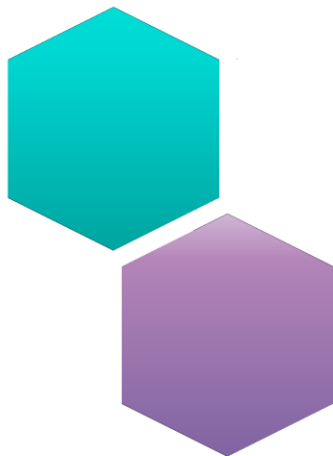
- Siemens collects 2EB data per day[7]
  - Sensor data can reveal corporate secrets (e.g., machine settings)
  - Predictive maintenance models usually trained on all data[8]

# Centralized Learning Approach



- does not scale with number of devices and neglects their computing power

# Centralized Learning Approach



- does not scale with number of devices and neglects their computing power
- prohibitive communication costs



# Centralized Learning Approach



- does not scale with number of devices and neglects their computing power
- prohibitive communication costs
- requires sharing of privacy-sensitive data

Chapter

# Federated Deep Learning

# Federated Deep Learning Approach



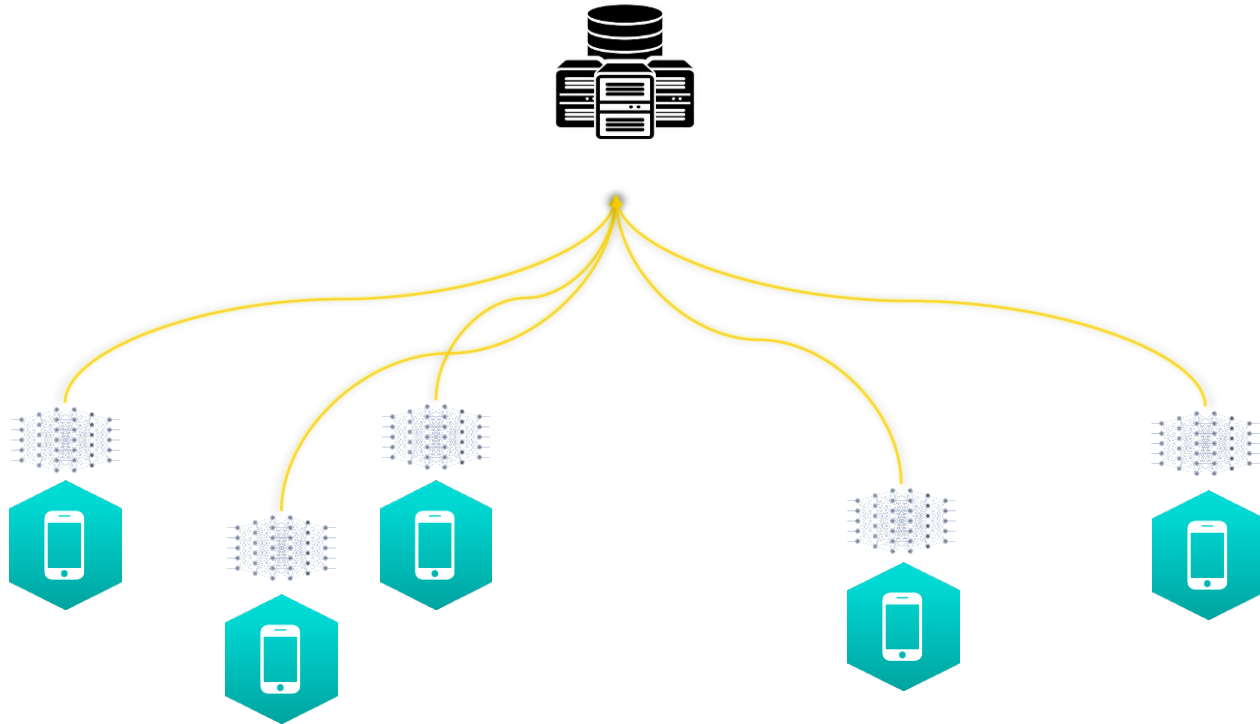
# Federated Deep Learning Approach



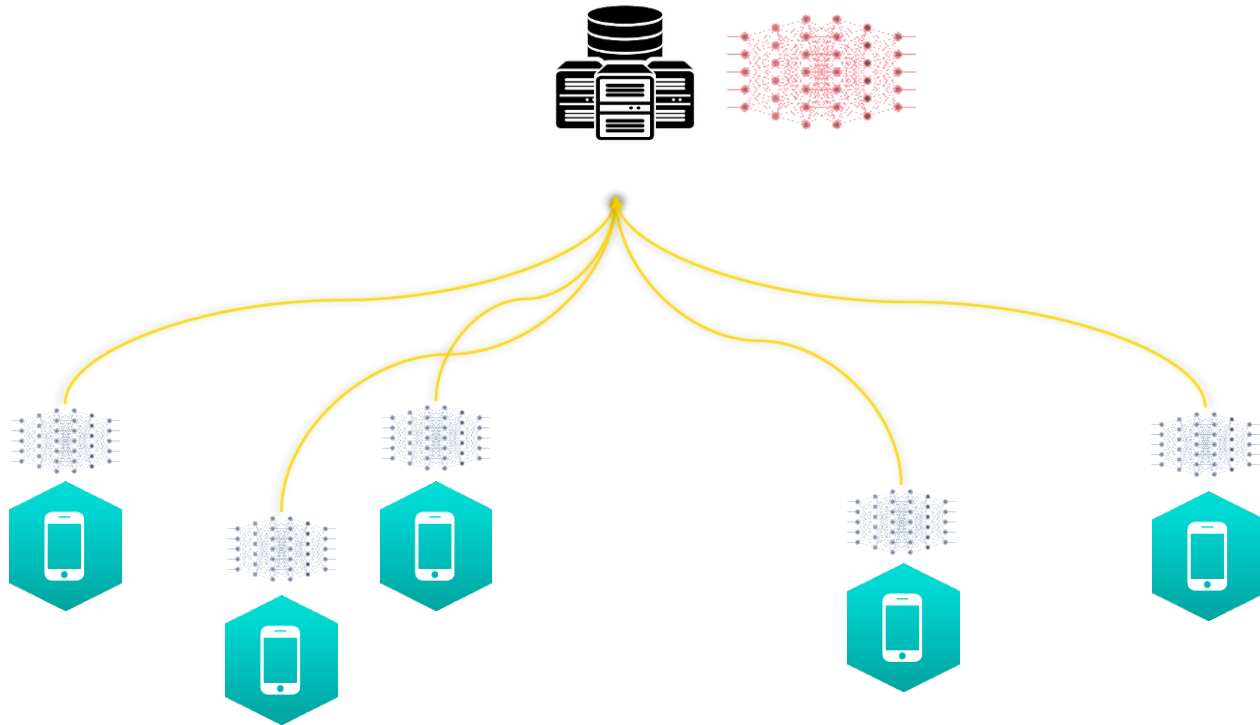
# Federated Deep Learning Approach



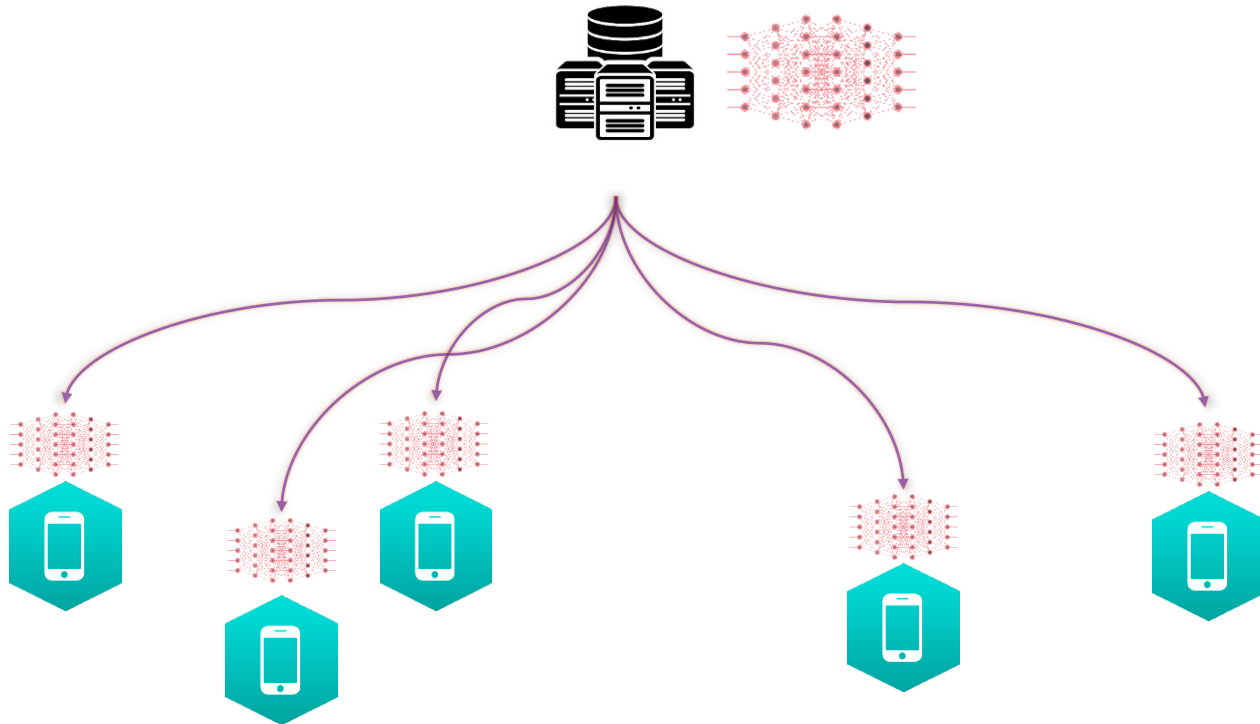
# Federated Deep Learning Approach



# Federated Deep Learning Approach

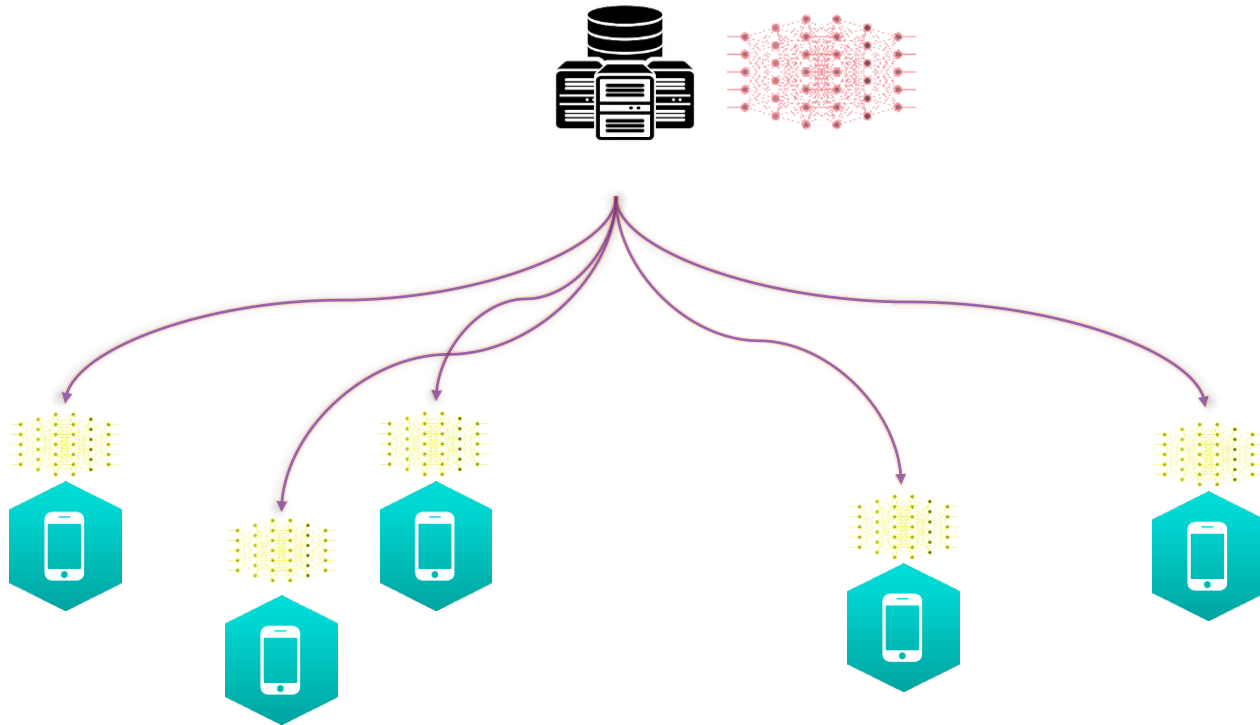


# Federated Deep Learning Approach





# Federated Deep Learning Approach



# Model Aggregation

$$\bar{f} = \frac{1}{m} \sum_{i=1}^m f^i$$

Average model

Number of learners

Local models

The diagram shows the mathematical formula for model aggregation:  $\bar{f} = \frac{1}{m} \sum_{i=1}^m f^i$ . Three green arrows point from text labels to parts of the formula: one from 'Average model' to  $\bar{f}$ , one from 'Number of learners' to the superscript  $m$  in the summation, and one from 'Local models' to the term  $f^i$ .

# Properties



- independent of the learning algorithm (black-box)

# Properties



- independent of the learning algorithm (black-box)
- works for large amount of model classes (linear + kernel models, various DNNs)

# Properties



- independent of the learning algorithm (black-box)
- works for large amount of model classes (linear + kernel models, various DNNs)
- requires that models can be trained locally, e.g., sufficient local computing power, and labels available

# Requirements

1. The Model trained in parallel has a high quality similar to one obtained by a the serial application of the learning algorithm.

# Requirements

1. The Model trained in parallel has a high quality similar to one obtained by a the serial application of the learning algorithm.
2. The parallel algorithm achieves a high speedup; ideally it has a polylogarithmic runtime on (quasi-) polynomial many learners.

# Requirements

1. The Model trained in parallel has a high quality similar to one obtained by a the serial application of the learning algorithm.
2. The parallel algorithm achieves a high speedup; ideally it has a polylogarithmic runtime on (quasi-) polynomial many learners.
3. The parallel execution requires communication that is almost linear in the number of learners and adaptive to the hardness of the learning problem.



# Requirements

1. The Model trained in parallel has a high **quality** similar to one obtained by a the serial application of the learning algorithm.
2. The parallel algorithm achieves a high **speedup**; ideally it has a polylogarithmic runtime on (quasi-) polynomial many learners.
3. The parallel execution requires **communication** that is almost linear in the number of learners and adaptive to the hardness of the learning problem.

# Privacy

- Additional inherent Feature
  - No Data is Shared between Parties
  - Only Model Weights are Shared
- Can be proven to be *differential private*[9]
  - By clipping the Gradient for each Weight

Chapter

**Leveraging Uncertainty in  
Federated Deep Learning**

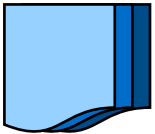
# Notation



- Communicator: Averages Models



- Worker: Trains Models



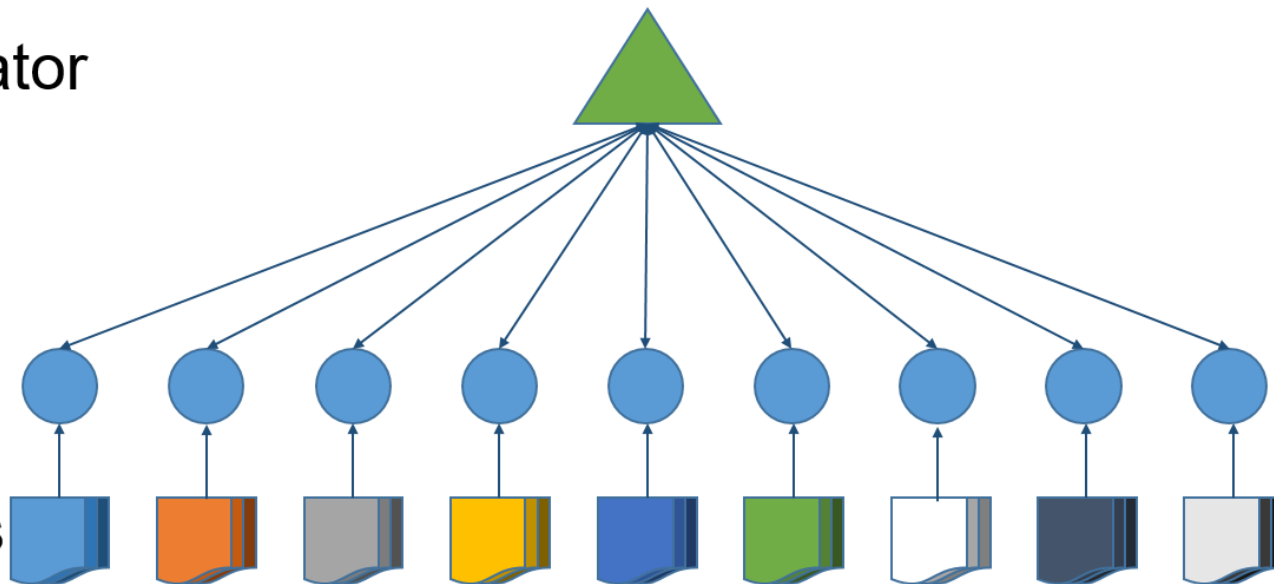
- Dataset: Contains Private Data

# Global Model (Baseline)

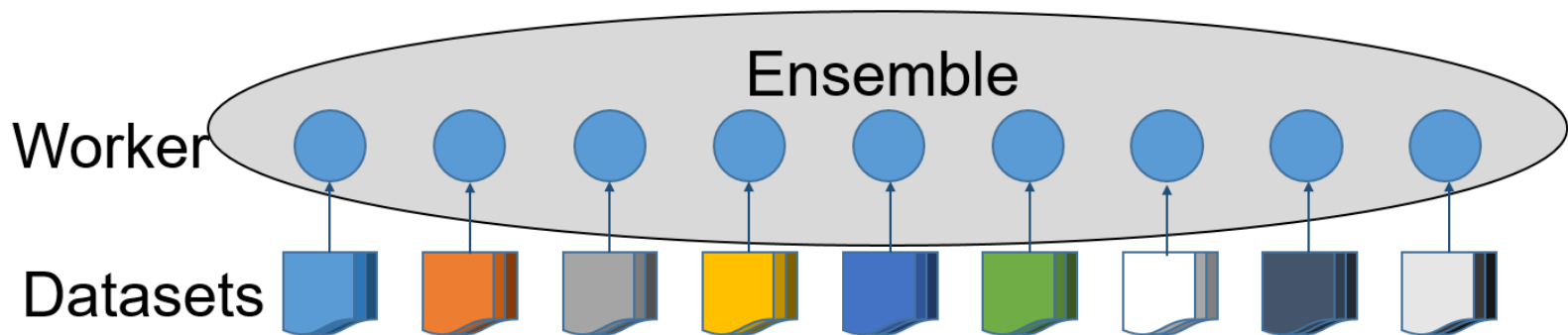
Coordinator

Workers

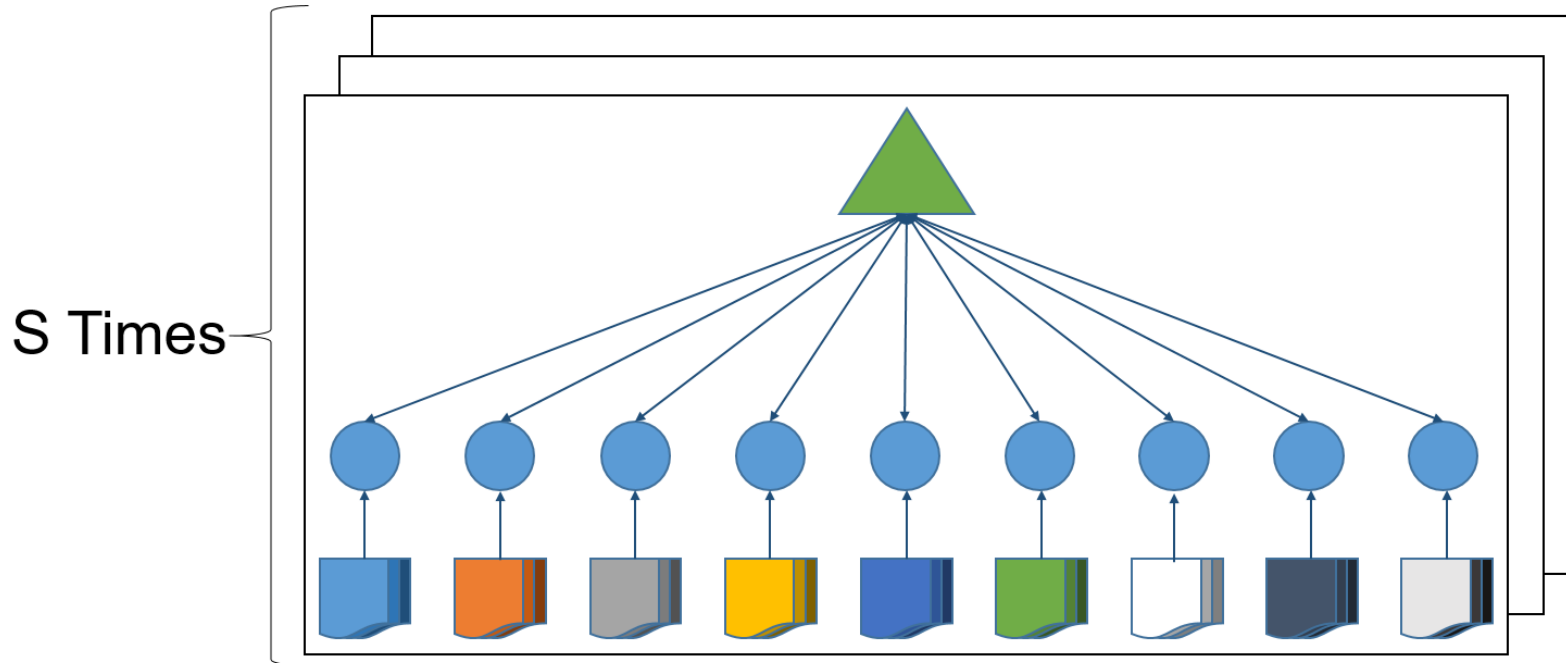
Datasets



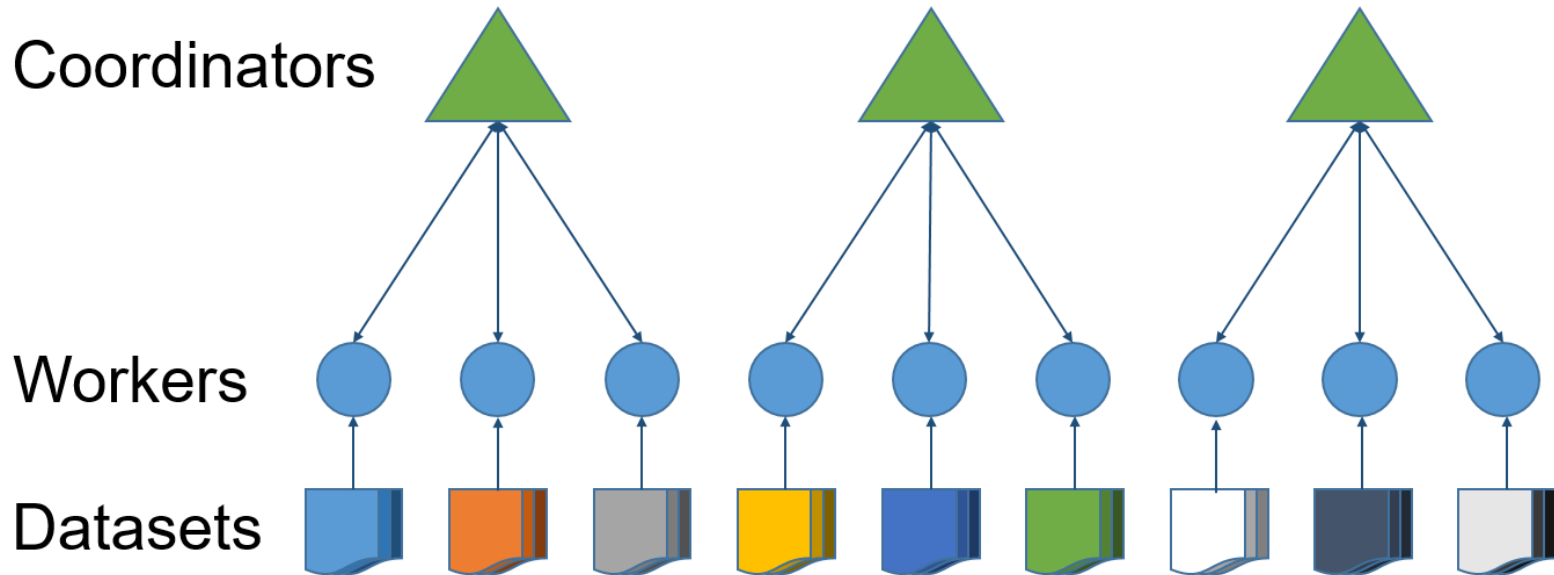
# Ensemble of Local Models



# Ensemble of Global Models



# Ensemble w. mult. Coord. Fixed Assignment





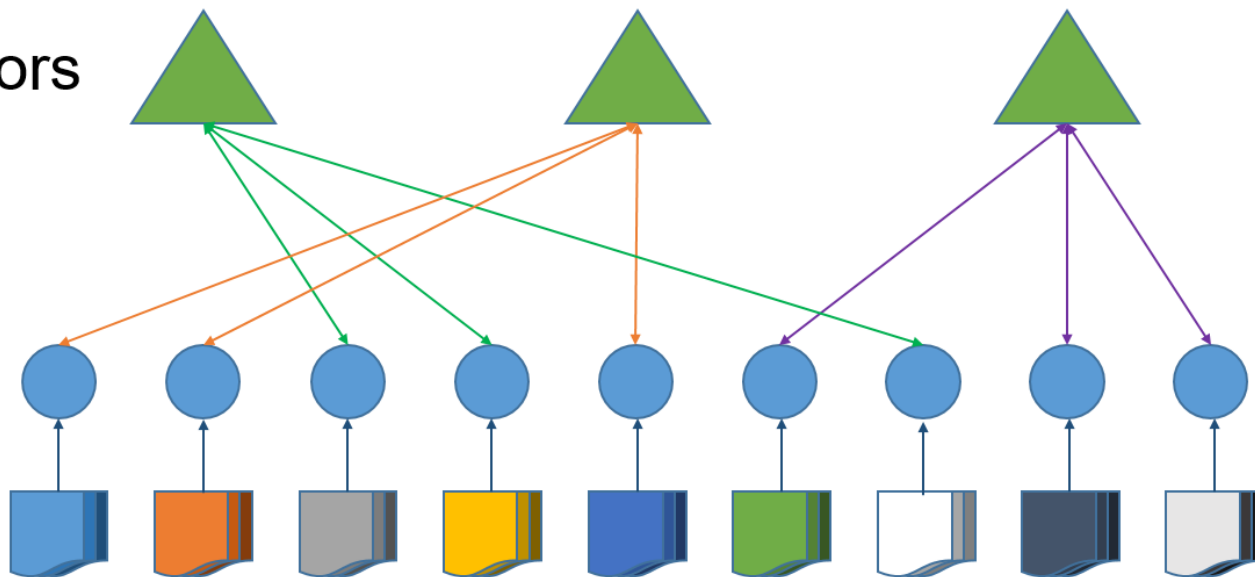
# Ensemble w. mult. Coord. Random Assignment

Coordinators

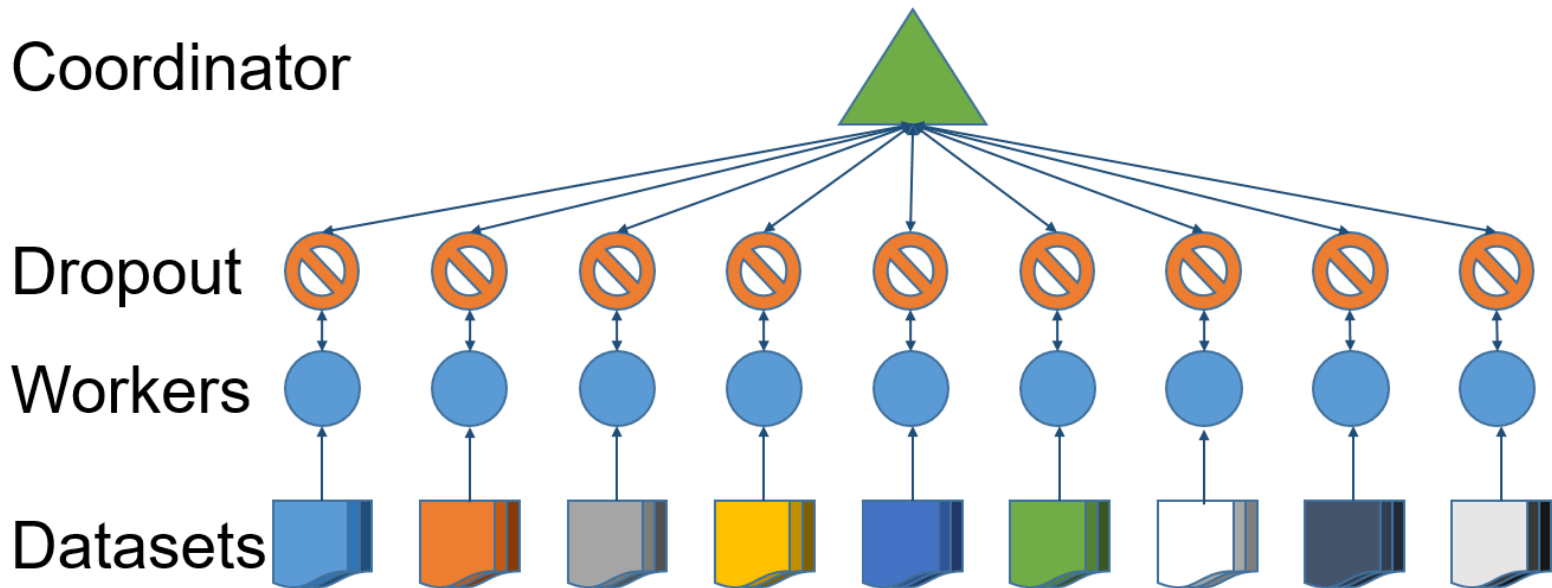


Workers

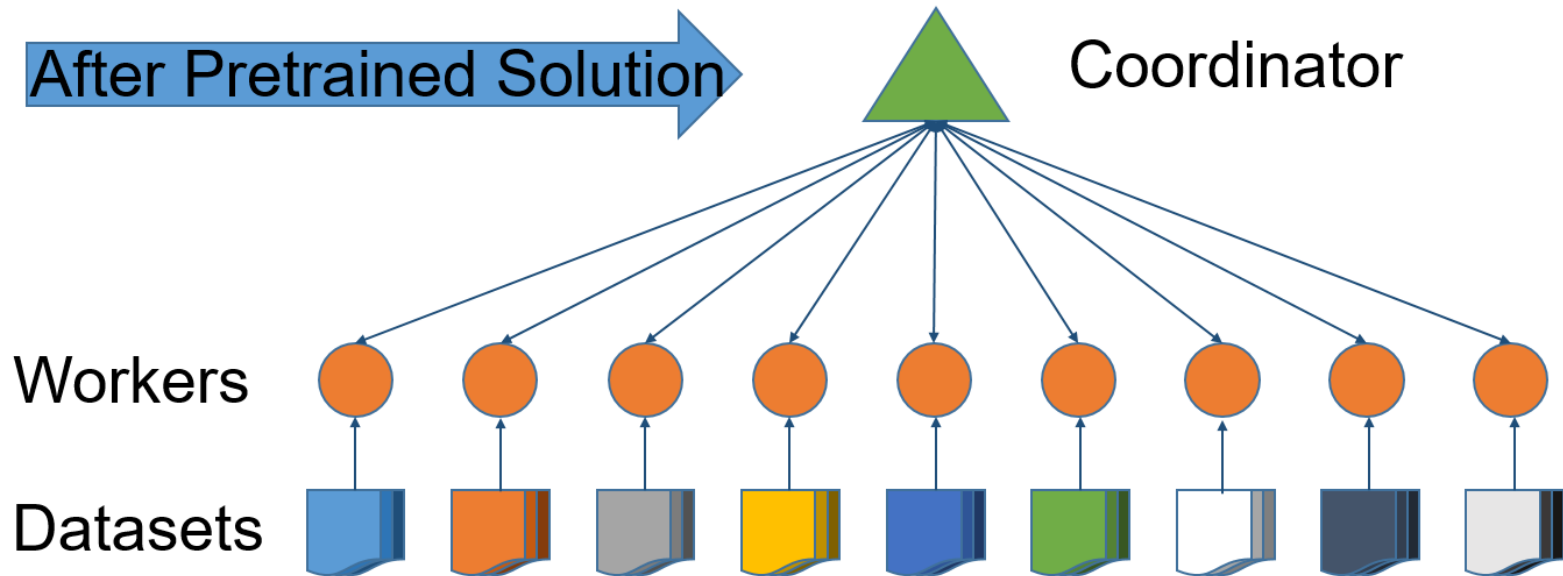
Datasets



# Federated MC-Dropout



# Federated SWAG



Chapter

**Empirical Evaluation**

# Empirical Evaluation

## Performance of Ensemble Approaches with CIFAR-10

	Detection of OOD (SVHN)			Detection of wrong classified Data	
	Accuracy	Entropy AUROC	Variance AUROC	Entropy AUROC	Variance AUROC
<b>Global Model</b>	86.58	0.924	-	0.888	-
<b>Ensemble of local models</b>	72.65	0.679	0.689	0.788	0.687
<b>Ensemble of global models</b>	<b>89.00</b>	<b>0.937</b>	<b>0.804</b>	<b>0.891</b>	<b>0.868</b>
<b>Ensemble w. mult. Coord.-FA</b>	83.43	0.860	0.740	0.858	0.815
<b>Ensemble w. mult. Coord.-RA</b>	86.72	0.920	0.750	0.887	0.861

# Empirical Evaluation

## Performance of MC-Dropout with CIFAR-10

	Accuracy	Detection of OOD (SVHN)		Detection of wrong classified Data	
		Entropy AUROC	Variance AUROC	Entropy AUROC	Variance AUROC
Global Model	77.88	0.742	-	0.838	-
Global Model with dropout	<b>86.58</b>	<b>0.924</b>	-	<b>0.888</b>	-
Federated MC-dropout	86.32	0.913	<b>0.714</b>	0.880	<b>0.849</b>

# Empirical Evaluation

## Performance of SWAG with CIFAR-10

	Accuracy	Detection of OOD (SVHN)		Detection of wrong classified Data	
		Entropy AUROC	Variance AUROC	Entropy AUROC	Variance AUROC
<b>Global Model</b>	86.58	0.924	-	0.888	-
<b>Federated SWAG</b>	<b>87.14</b>	<b>0.924</b>	<b>0.807</b>	<b>0.892</b>	<b>0.856</b>

# Chapter

# **Conclusion**



# Conclusion

- Best Results in terms of accuracy and AUROC scores
  - Ensemble of Local Models
    - Best Accuracy
    - High Computational Overhead
  - Federated SWAG
    - Easy Integration in Running System
    - Doubled Communication during SWAG Training

# Conclusion

- Multiple Coordinators
  - No Performance Increase
  - Might be suitable for moving devices like autonomous driving
- Dropout
  - Improves Performance when used as regularization
  - MC-Dropout does not improve performance

# Chapter Discussion

[florian.linsner@rub.de](mailto:florian.linsner@rub.de)