

Structured Cooperative Learning with Graphical Model Priors

Shuangtong Li¹, Tianyi Zhou², Xinmei Tian^{1 3}, Dacheng Tao⁴

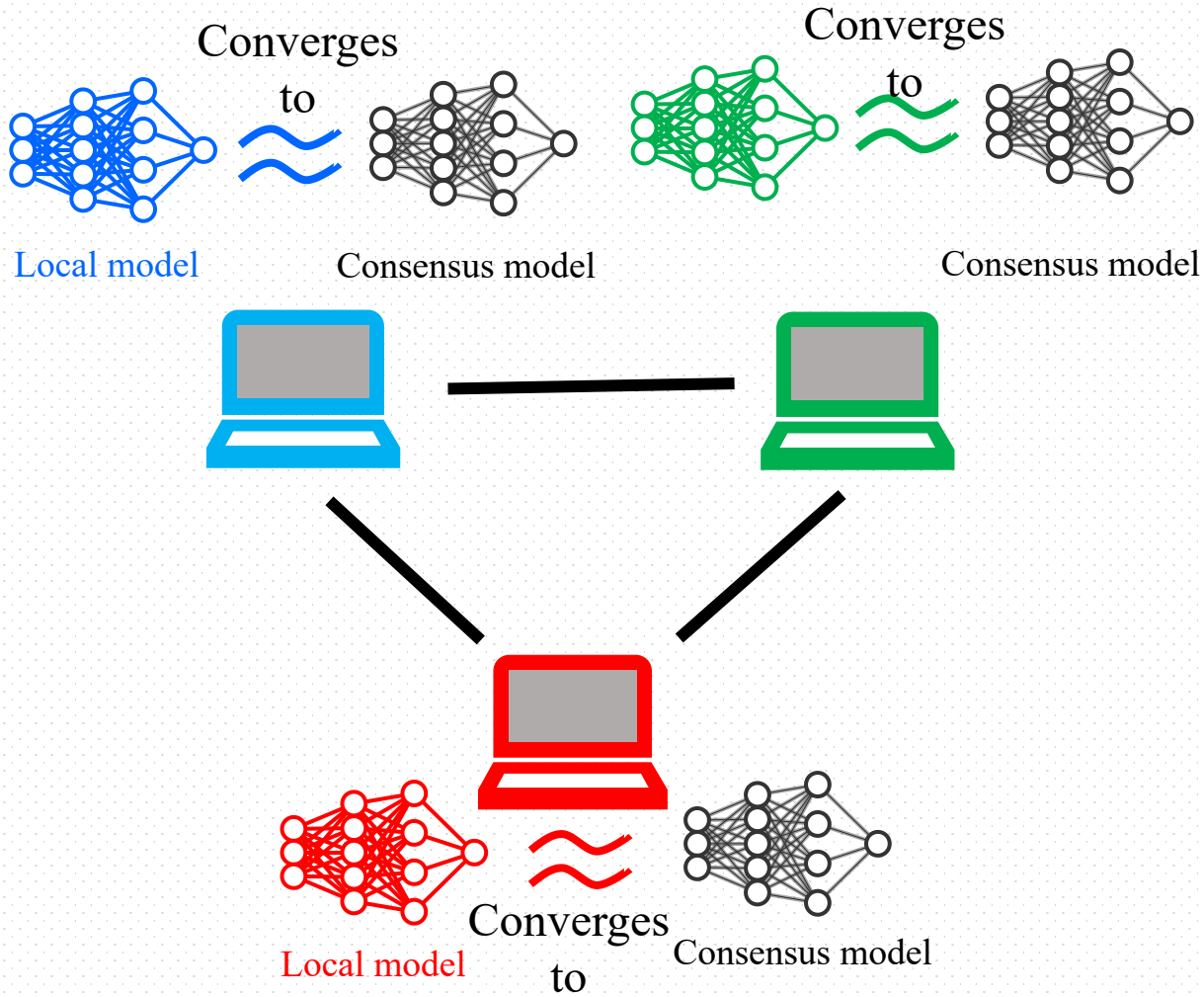
¹University of Science and Technology of China

²University of Maryland, College Park

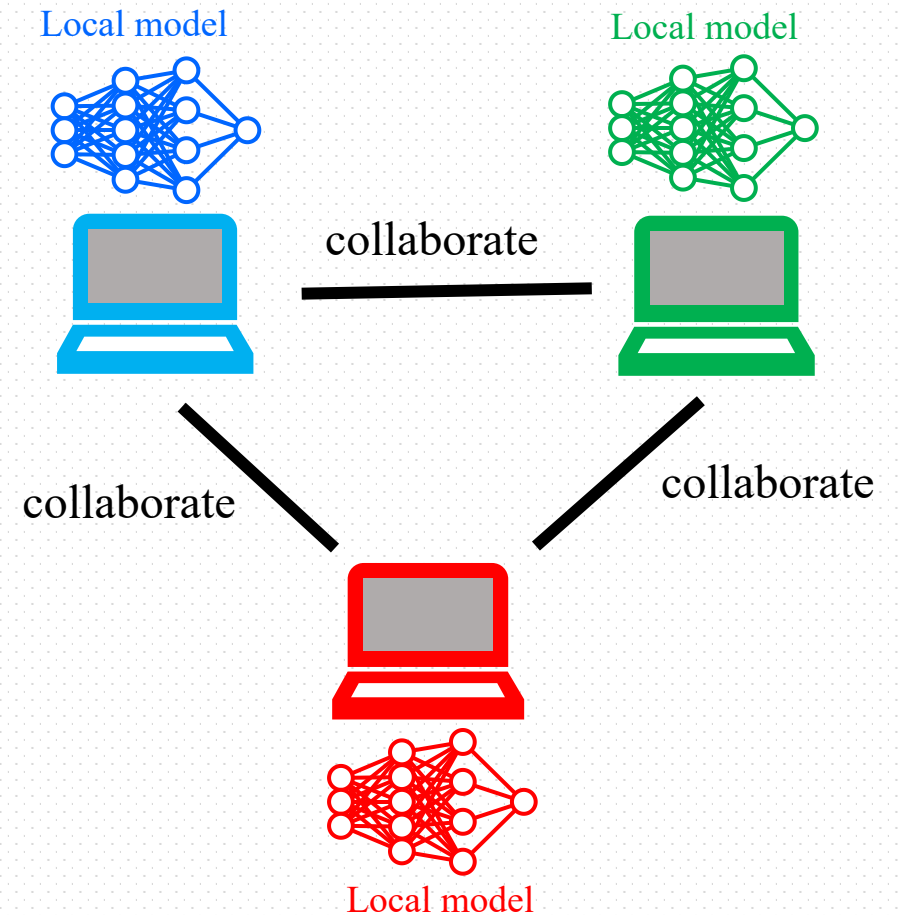
³Institute of Artificial Intelligence, Hefei Comprehensive National Science Center

⁴The University of Sydney

Background



Traditional decentralized Learning



Decentralized Learning of Personalized Models (DLPM) [1]

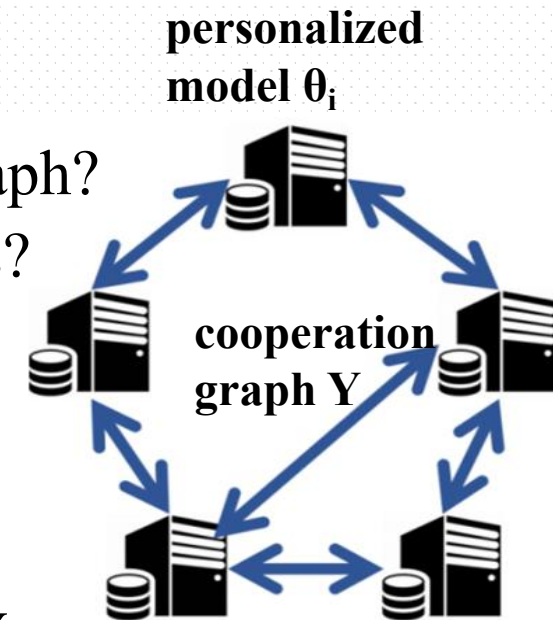
Motivation

DLPM Challenges:

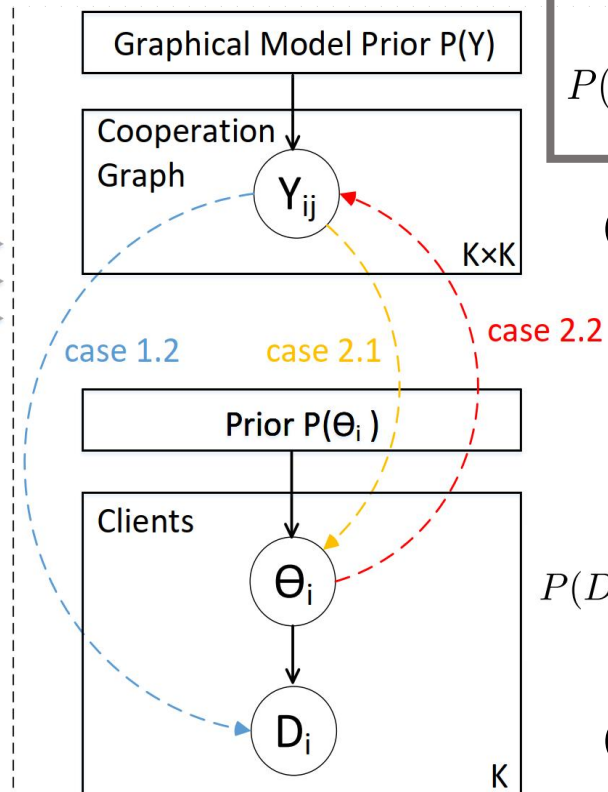
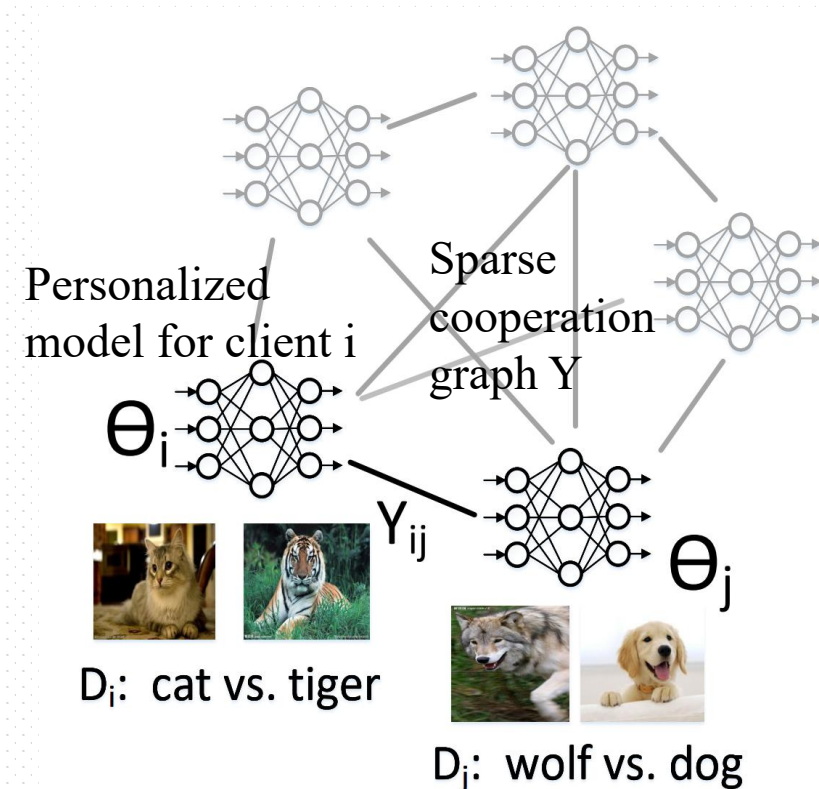
- How to determine when and which clients should cooperate?
- How to cooperate when personal tasks and data cannot be shared?
- To save communication cost, how to discover a sparse cooperation graph?
- How to adjust the graph adaptive to model changes in training process?

Structured Cooperative Learning (SCoolL):

- ❖ A **general** probabilistic modelling framework.
- ❖ Jointly optimize **personalized models $\theta_{1:K}$** and **cooperation graph Y** .
- ❖ Different graphical model priors of Y \longrightarrow various novel DLPM algorithms.
- ❖ A systematic optimization method: **variational inference**.



SCoolL Framework



Probabilistic Modeling with Cooperation Graph

$$P(\theta_{1:K} | D_{1:K}) \propto P(\theta_{1:K}, D_{1:K}) = \int \boxed{P(D_{1:K} | \theta_{1:K}, Y)} \boxed{P(\theta_{1:K}, Y)} dY.$$

(1) Joint Likelihood $\boxed{P(D_{1:K} | \theta_{1:K}, Y)}$

case 1.1 Y does not affect data distribution.

$$P(D_{1:K} | \theta_{1:K}) = \prod_{i=1}^K P(D_i | \theta_i)$$

case 1.2 Y coordinates the training process.

$$P(D_{1:K} | \theta_{1:K}, Y) = \prod_{i=1}^K P(D_{1:K} | \theta_i, Y) = \prod_{i=1}^k \left(P(D_i | \theta_i) \prod_{j \neq i, Y_{ij}=1} P(D_j | \theta_i) \right)$$

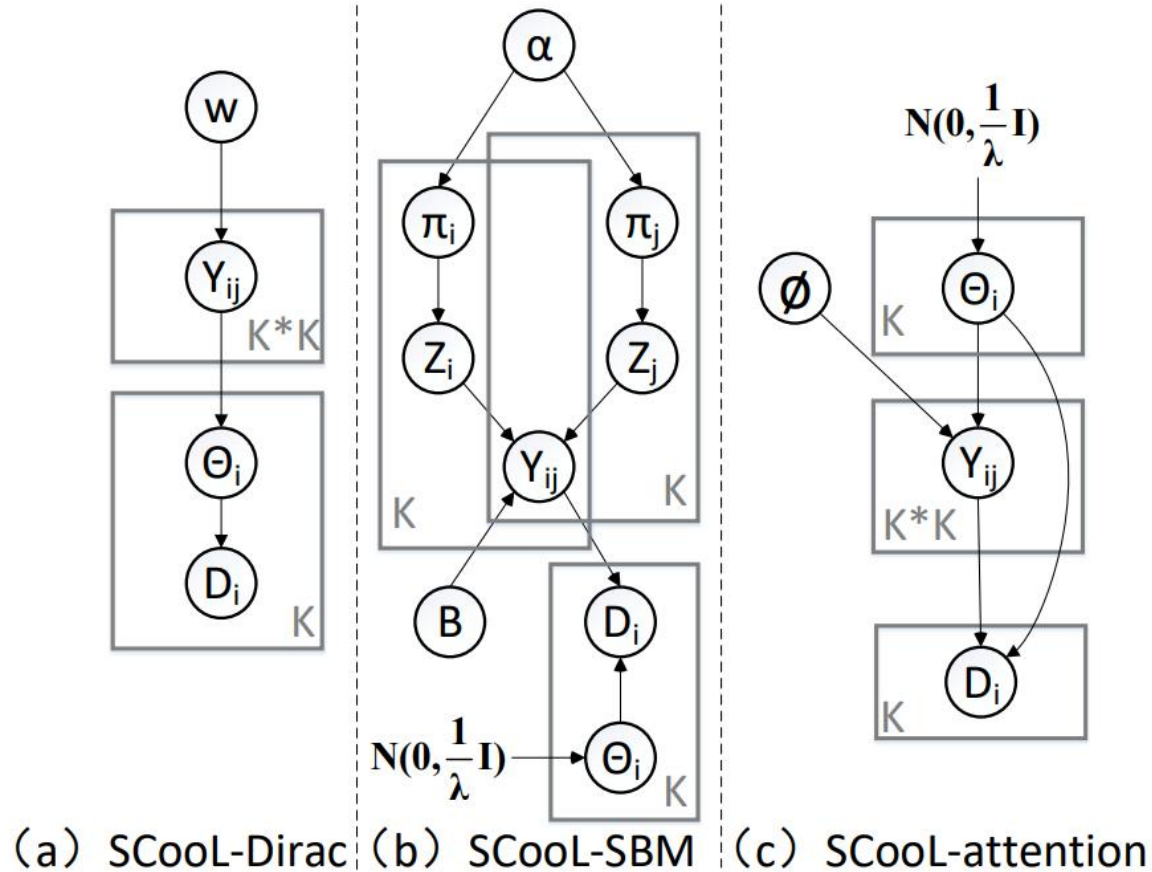
(2) Joint Priors $\boxed{P(\theta_{1:K}, Y)}$

case 2.1 : $P(\theta_{1:K} | Y)P(Y)$, $\theta_{1:K}$ is derived from Y .

case 2.2 : $P(Y | \theta_{1:K})P(\theta_{1:K})$, $\theta_{1:K}$ determines Y .

case 2.3 : $P(\theta_{1:K})P(Y)$, $\theta_{1:K}$ is independent to Y .

SCool Instantiations



(a) SCool-Dirac

$$Y \sim \delta(w)$$

SCool-Dirac is equivalent to DPSGD [2].

(b) SCool-SBM

$Y \sim$ Stochastic Block model (SBM) [3]

(c) SCool-attention

$$Y_{ij} \sim \text{Attention}(\theta_i, \theta_j)$$

$$p_{ij} = \frac{\exp(f(\theta_i, \theta_j))}{\sum_l \exp(f(\theta_i, \theta_l))}$$

$$\vec{Y} \sim \text{Categorical}(p_{i1}, \dots, p_{ik})$$

[2] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In Advances in Neural Information Processing Systems, 2017.

[3] Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. Social networks, 5(2):109–137, 1983.

EM Algorithm for SCool

We derive EM algorithms for SCool models via **variational inference** method.

ELBO:

$$\begin{aligned}\log p(X|\Phi) &= \log \int p(X, Z|\Phi) dZ \\ &\geq \int q(Z) \log \frac{p(X, Z|\Phi)}{q(Z)} dZ := H(q, \Phi).\end{aligned}$$

E-step: update cooperation graph Y .

$$w_{ij} \leftarrow F\left(\log P(D_j|\theta_i), \beta, \Phi\right) \forall i, j \in [K]$$

M-step: optimize the local models $\theta_{1:K}$.

$$\theta_i \leftarrow \theta_i - \eta_1 \left(\sum_{j \neq i} w_{ij} \nabla L(D_j; \theta_i) + \nabla L(D_i; \theta_i) + G(\beta, \Phi) \right)$$

Experiments

Methodology	Algorithm	CIFAR-10	CIFAR-100	MinImageNet
Local only	local SGD	87.5 \pm 7.02	55.47 \pm 5.20	41.59 \pm 7.71
Federated	FedAvg	70.65 \pm 10.64	40.15 \pm 7.25	34.26 \pm 6.01
	FOMO	88.72 \pm 5.41	52.44 \pm 5.09	44.56 \pm 4.31
	Ditto	87.32 \pm 6.42	54.28 \pm 5.31	42.73 \pm 5.19
Decentralized	D-PSGD(1s)	83.01 \pm 7.34	40.56 \pm 6.94	30.26 \pm 5.75
	D-PSGD(5e)	75.89 \pm 6.65	35.03 \pm 4.83	28.41 \pm 5.18
	CGA(1s)	65.65 \pm 12.66	30.81 \pm 10.79	27.65 \pm 11.78
	CGA(5e)	diverge	diverge	diverge
	SPDB(1s)	82.36 \pm 7.14	54.29 \pm 6.15	39.17 \pm 3.93
	SPDB(5e)	81.15 \pm 7.06	53.23 \pm 7.48	35.93 \pm 5.05
	Dada	85.65 \pm 6.36	57.61 \pm 5.45	37.81 \pm 7.15
	meta-L2C	92.10 \pm 4.71	58.28 \pm 3.09	48.80 \pm 4.17
SCool (Ours)	SCool -SBM	91.37 \pm 5.03	58.76 \pm 4.30	48.69 \pm 5.21
	SCool -attention	92.21\pm5.15	59.47\pm4.95	49.53\pm3.29

Conclusion

we propose a **general probabilistic modelling framework**: Structured Cooperative Learning (SCool), for DLPM problems.

- ❖ SCool jointly optimizes **personalized models $\theta_{1:K}$** and **cooperation graph Y** .
- ❖ Different graphical model priors of Y generate various novel DLPM algorithms.
- ❖ SCool uses a systematic optimization method: **variational inference**.
- ❖ SCool outperforms previous federated / decentralized learning baselines in experiments.