

# Distributed Weight Consolidation: A Brain Segmentation Case Study

Patrick McClure<sup>1</sup>, Charles Y. Zheng<sup>1</sup>, Jakub R. Kaczmarzyk<sup>2</sup>, John A. Lee<sup>1</sup>, Satrajit S. Ghosh<sup>2</sup>, Dylan Nielson<sup>1</sup>, Peter Bandettini<sup>3</sup>, Francisco Pereira<sup>1</sup>

1) Machine Learning Team, National Institute of Mental Health; 2) McGovern Institute for Brain Research, Massachusetts Institute of Technology; 3) Section on Functional Imaging Methods, National Institute of Mental Health

Summary
<b>Problem:</b> Training a network on distributed data from multiple sites that cannot share their data
<b>Potential Solution:</b> Use continual learning to sequentially update a network one site at a time using the previous network as a prior
<b>Drawback:</b> Each site will have to wait for all previous sites before training a network on their data
<b>Our Proposal:</b> Use continual learning at each site and then combine the resulting networks into a new prior for further training

Methods
<b>Variational Bayesian Neural Networks [6,7,8,9]</b>
$p(W D_1) = \frac{p(D_1 W)p(W)}{p(D_1)}$
$\operatorname{argmin}_V - \int q_V(W) \log p(D_1 W) dW + KL[q_V(W)  p(W)]$
<b>Variational Continual Learning (VCL) [10]</b>
$p(W D_{1:T}) = \frac{p(D_T W)p(W D_{1:T-1})}{p(D_T)}$
$\operatorname{argmin}_V - \int q_V(W) \log p(D_T W) dW + KL[q_V(W)  p(W D_{1:T-1})]$
<b>Distributed Weight Consolidation (DWC)</b>
$D_T = \{D_T^1, \dots, D_T^S\}$
$p(W D_{1:T-1}, D_T^S) = \frac{p(D_T^S W)p(W D_{1:T-1})}{p(D_T^S)}$
$p(W D_{1:T}) = \frac{1}{p(W D_{1:T-1})^{S-1}} \prod_{s=1}^S p(W D_{1:T-1}, D_T^S)$
$\operatorname{argmin}_V - \int q_V(W) \log p(D_{T+1} W) dW + KL[q_V(W)  p(W D_{1:T})]$

Visual Results
<b>HCP</b>
<b>NKI</b>

Data and Architecture
Labels [1]
0 "background" 1 "white matter" ⋮ 10 "hippocampus" 11 "amygdala" ⋮ 48 "ctx-transversetemporal" 49 "ctx-insul"
<b>Data</b>
HCP [2]: ~900 training and ~100 test MRI volumes
NKI [3]: ~1000 training and ~100 test MRI volumes
Buckner [3]: ~200 training and ~20 test MRI volumes
WU120 [4]: ~100 training and ~10 test MRI volumes
ABIDE [5]: ~200 test MRI volumes
<b>Architecture</b>
MeshNet [1,5]: An 8-layer 3d dilated convolutional neural network

Numerical Results						
Average Dice Score across Classes						
Network	H	N	B	W	Avg.	A
$H_{MAP}$	82.25	65.88	67.94	70.88	72.92	55.25
$N_{MAP}$	71.2	72.19	70.73	73.06	71.66	66.67
$B_{MAP}$	65.69	50.17	82.02	68.87	59.25	50.23
$W_{MAP}$	70.18	66.27	72.2	76.38	59.25	62.83
$H \rightarrow N$	75.40	73.24	71.77	73.17	74.03	64.62
$H \rightarrow B$	73.85	56.79	79.49	68.53	65.78	49.27
$H \rightarrow W$	77.07	67.63	76.15	77.26	67.63	62.31
$H \rightarrow N \rightarrow B \rightarrow W$	77.42	71.46	79.70	79.82	74.86	63.3
$H \rightarrow W \rightarrow B \rightarrow N$	78.04	78.15	75.79	79.50	77.99	70.79
$H \rightarrow N + B + W$ (DWC w/o FT)	78.28	73.52	78.02	77.37	75.95	65.56
<b>Ensemble</b>	79.13	72.32	80.02	78.84	75.94	66.27
$H \rightarrow N + B + W \rightarrow H$ (DWC)	80.34	73.64	77.46	78.10	76.82	66.21
$HNBW_{MAP}$	81.38	77.99	80.64	79.54	79.62	70.76

Conclusions
There are many problems for which accumulating data into one accessible dataset for training can be difficult or impossible, such as for clinical data. It may, however, be feasible to share models derived from such data.
DWC:
<ul style="list-style-type: none"> <li>Allows for using Bayesian continual learning non-sequentially</li> <li>Performs as well or better than or better than the MAP ensemble</li> </ul>

## References