

A Theoretical Framework for Deep and Locally Connected ReLU Network

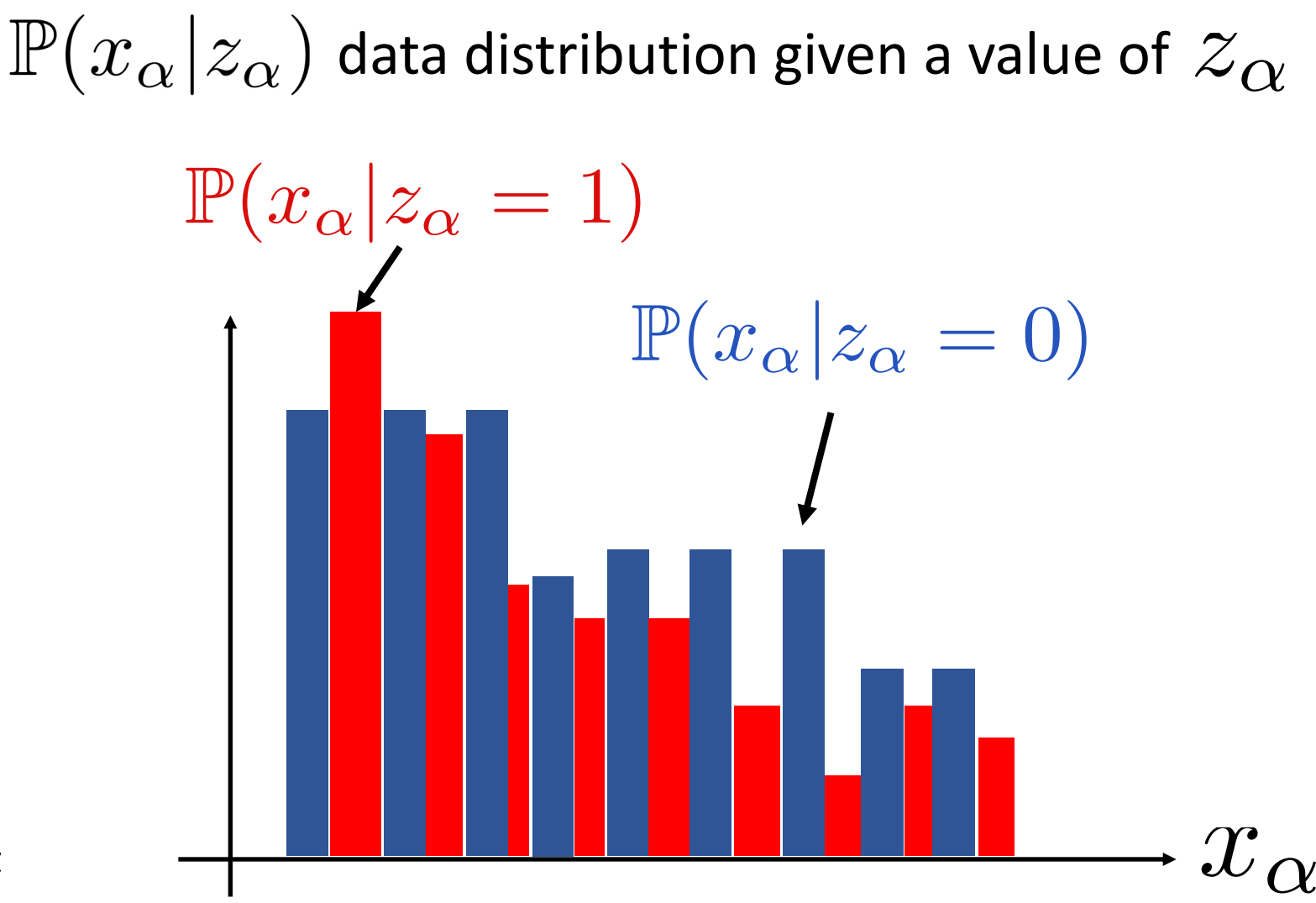
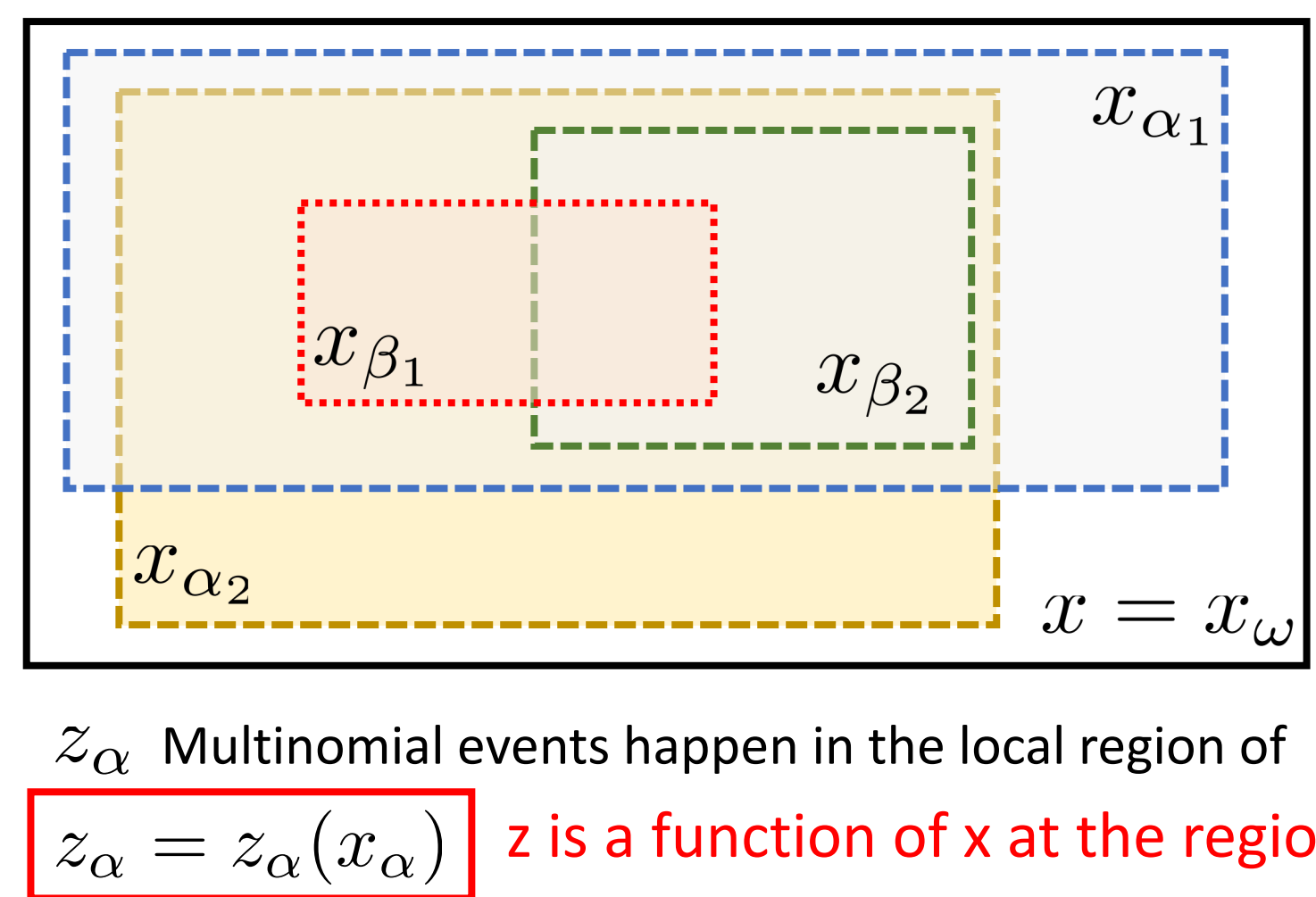
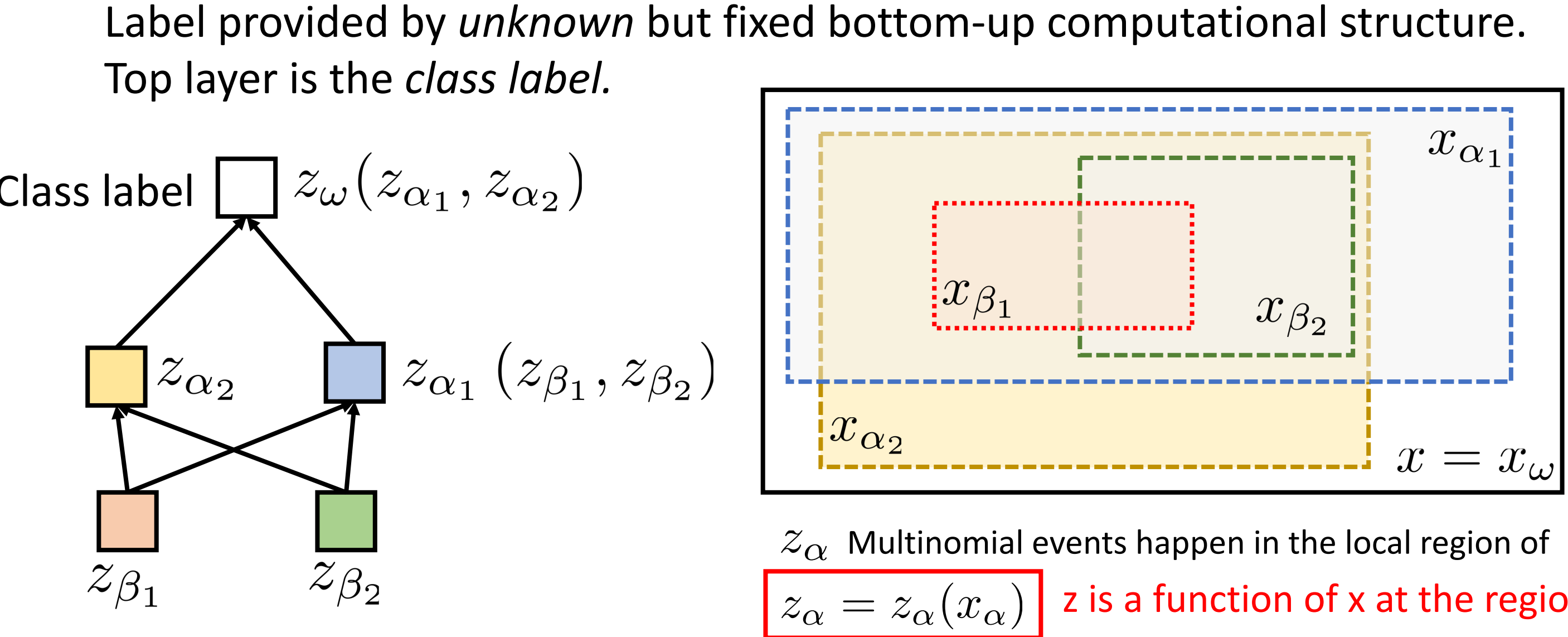
Yuandong Tian, Facebook AI Research

Arxiv: <https://arxiv.org/abs/1809.10829>

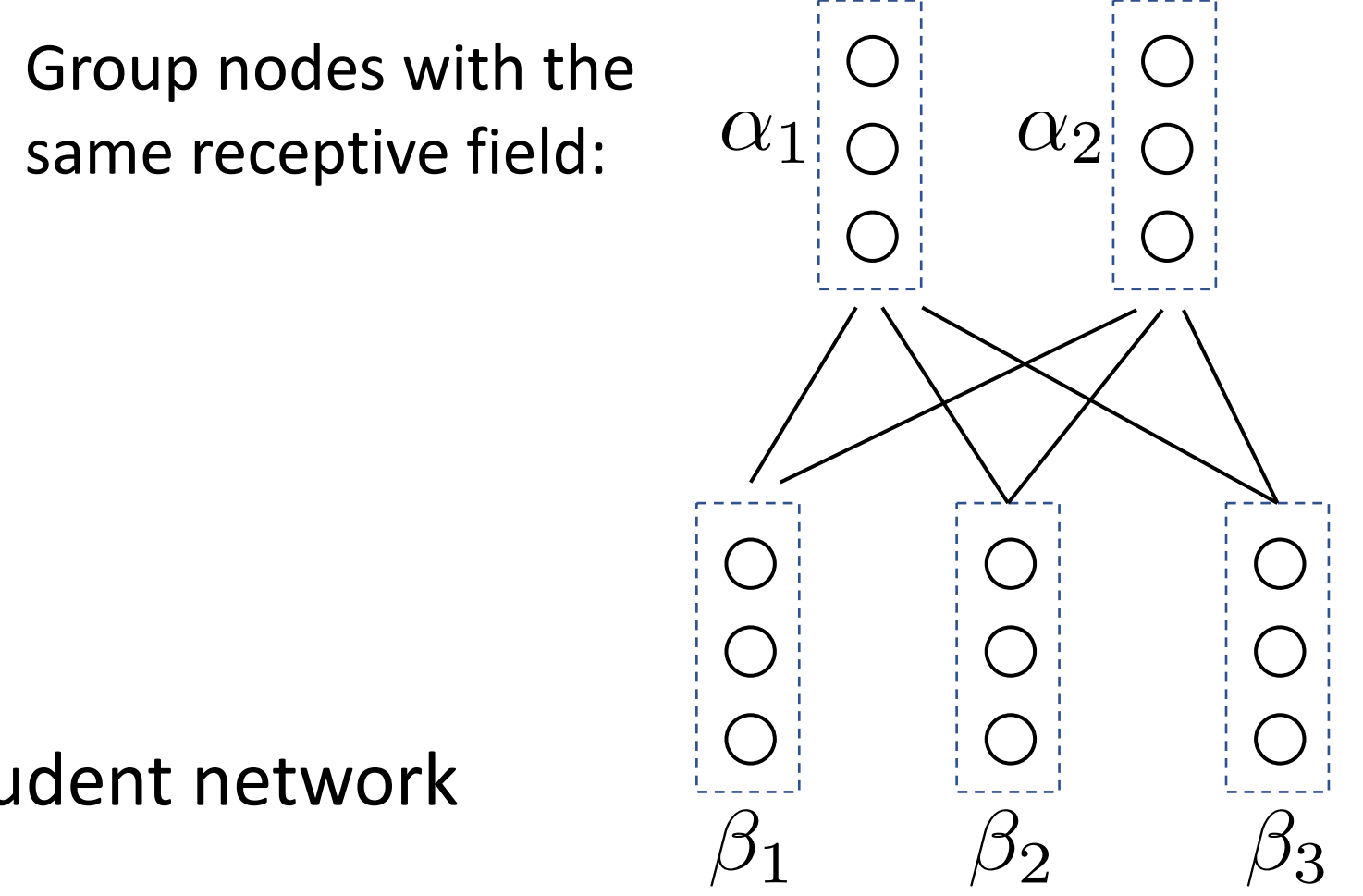
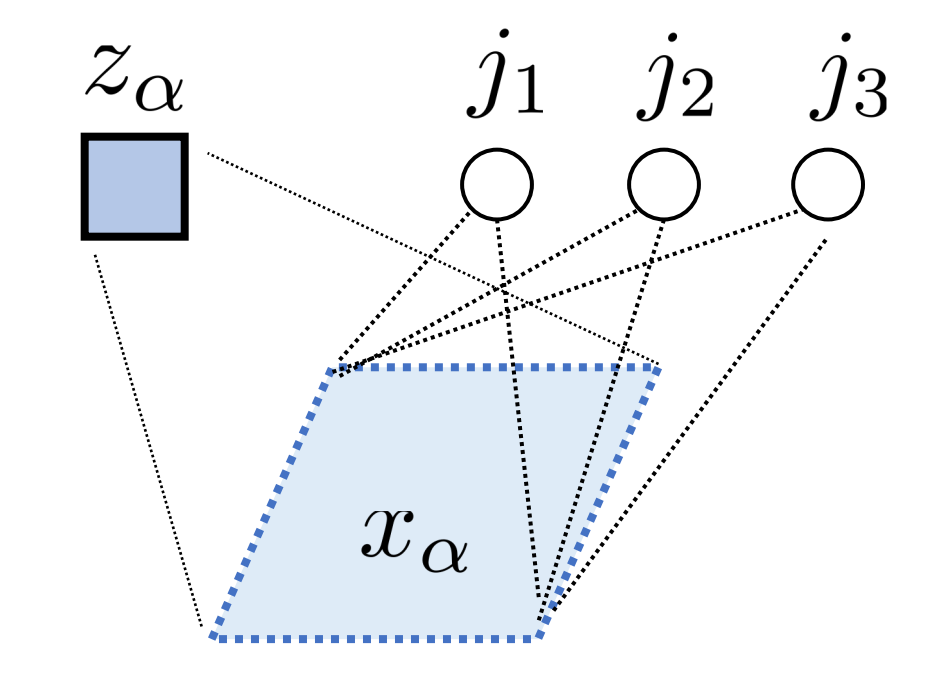
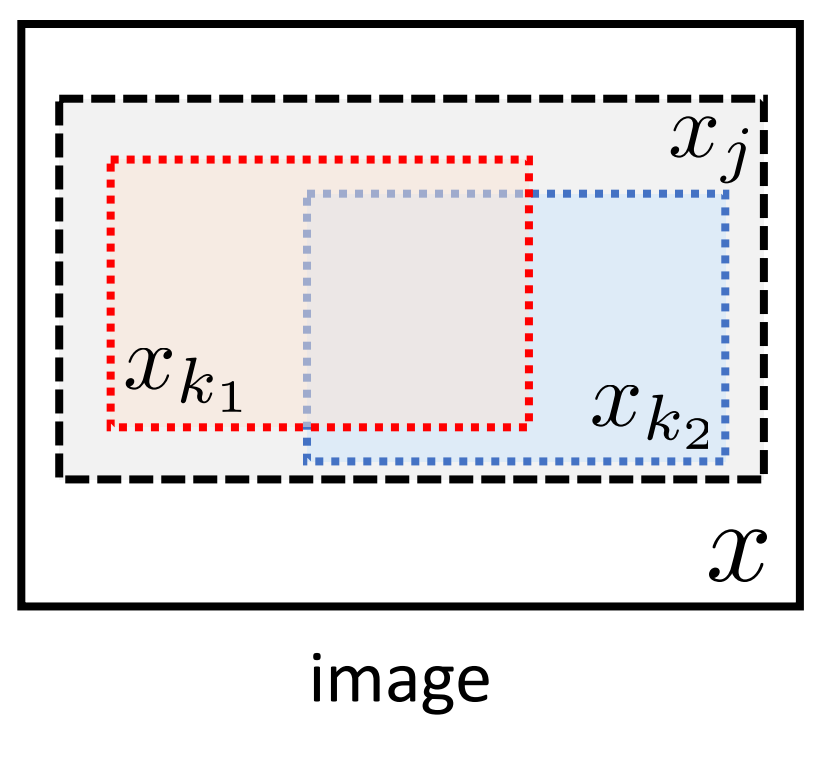
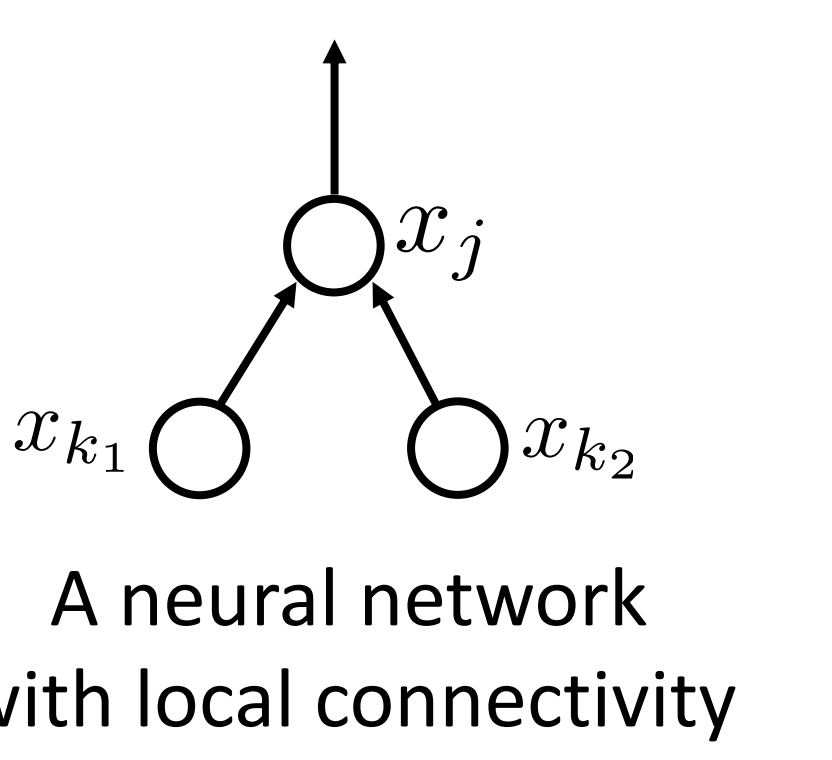


Problem Setting

Teacher (Locally connected)



Student Network



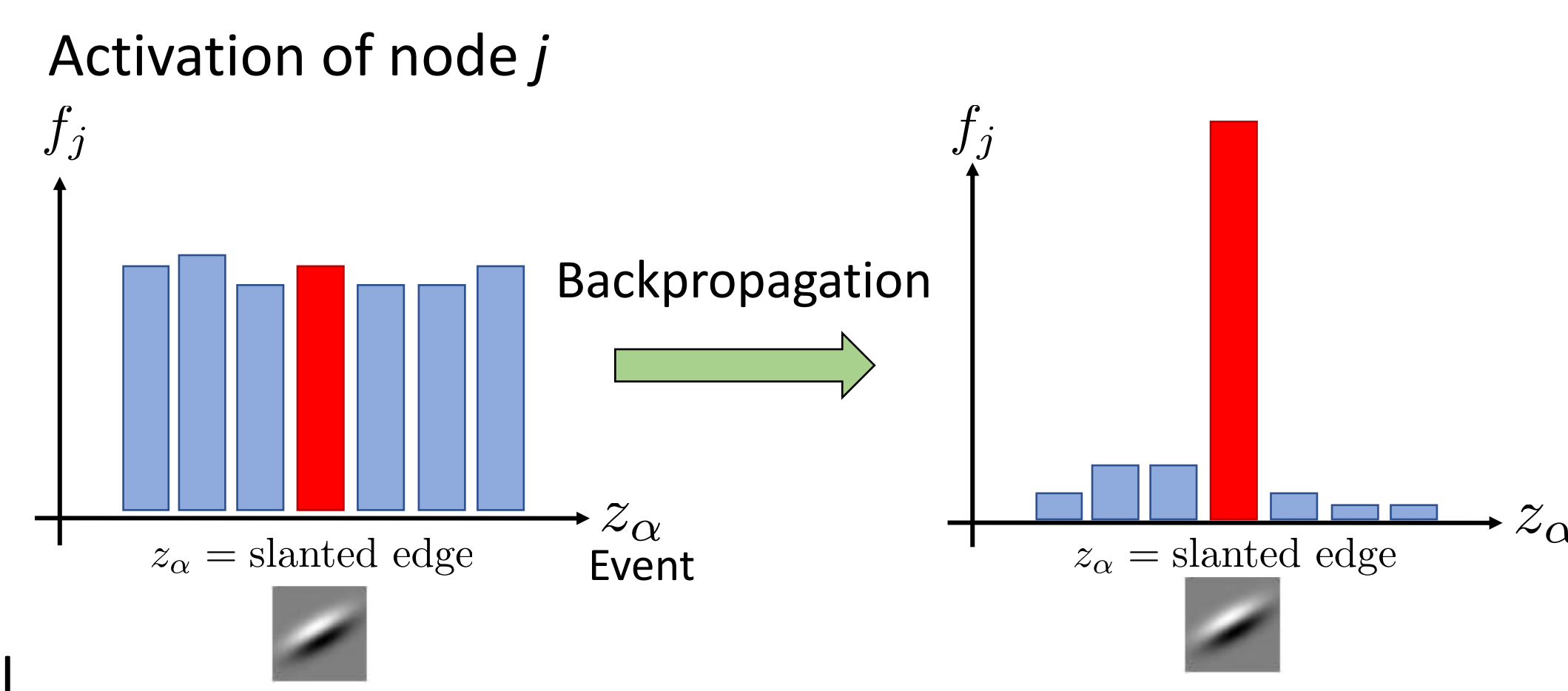
Objective: Match the output of the student network with class labels.

Goal

Student receives supervision from the highest level
Will student mimic the teacher **at every level** after GD?

Define: $f_j(z_\alpha) \equiv \mathbb{E}_{x|z_\alpha} [f_j(x)]$

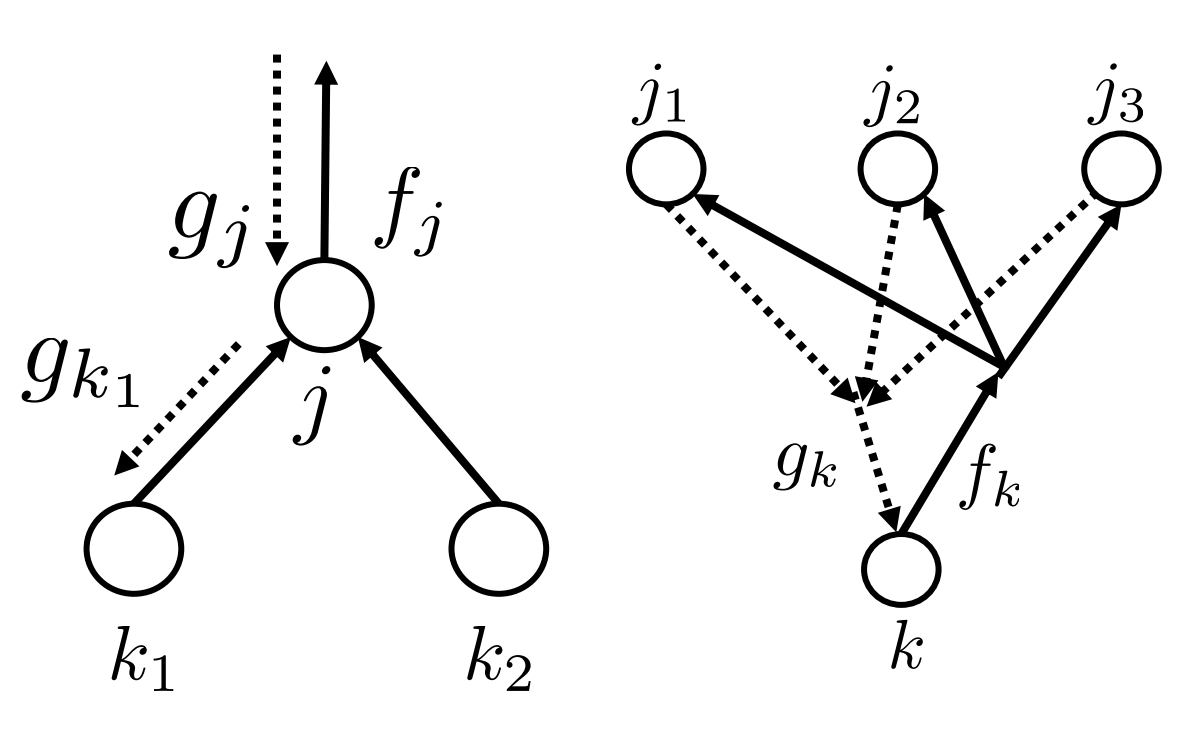
What we want: $f_j(z_\alpha = a) > 0$
 $f_j(z_\alpha \neq a) \approx 0$



Long term [Super hard]: Show propagation can achieve this goal.

First step [This paper]: Build a formulation first based on this setting and discuss about its properties

Forward and Backward Propagation in Student Network



Notation of scalar quantities (all dependent on weights)

$f_j(x)$ Activation of node j when input is x

$f'_j(x)$ Gating of node j when input is x

$g_j(x)$ Gradient propagated to node j , when input is x

W_{jk} Weights connecting node j to node k

$f_j(x) = f'_j(x) \sum_{k \in \text{ch}(j)} w_{jk} f_k(x)$

$g_k(x) = f'_k(x) \sum_{j \in \text{pa}(k)} w_{jk} g_j(x)$

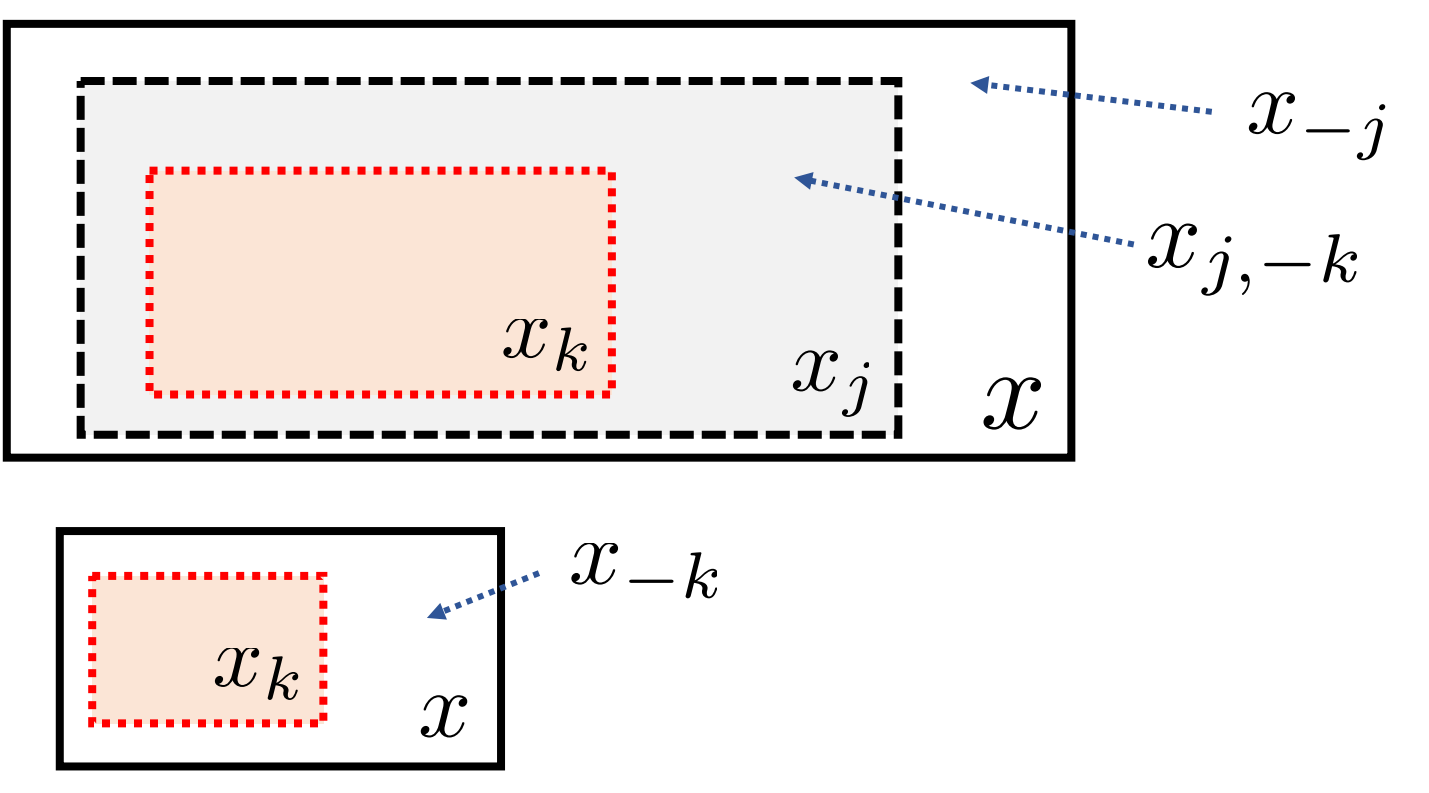
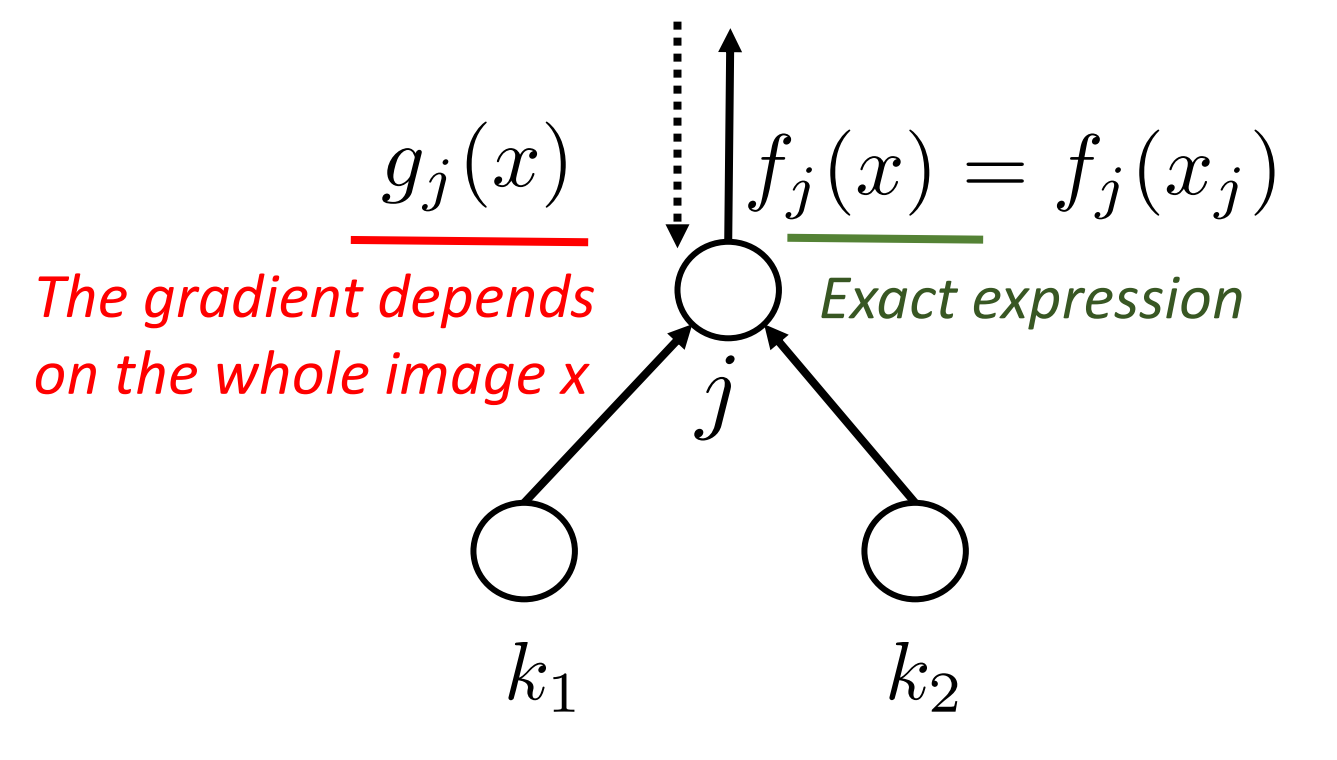
$\Delta w_{jk} = \mathbb{E}_x [f_k(x) g_j(x)]$

Using locality

Marginalized gradient

$g_j(x_k) = \mathbb{E}_{x_{-k}|x_k} [g_j(x)]$

$g_k(x_k) = \mathbb{E}_{x_{-k}|x_k} [g_k(x)]$



Theorem 1 (Recursive Property of marginalized gradient). $g_j(x_k) = \mathbb{E}_{x_{j,-k}|x_k} [g_j(x_j)]$

Reformulation

Theorem 2 [Reformulation]:

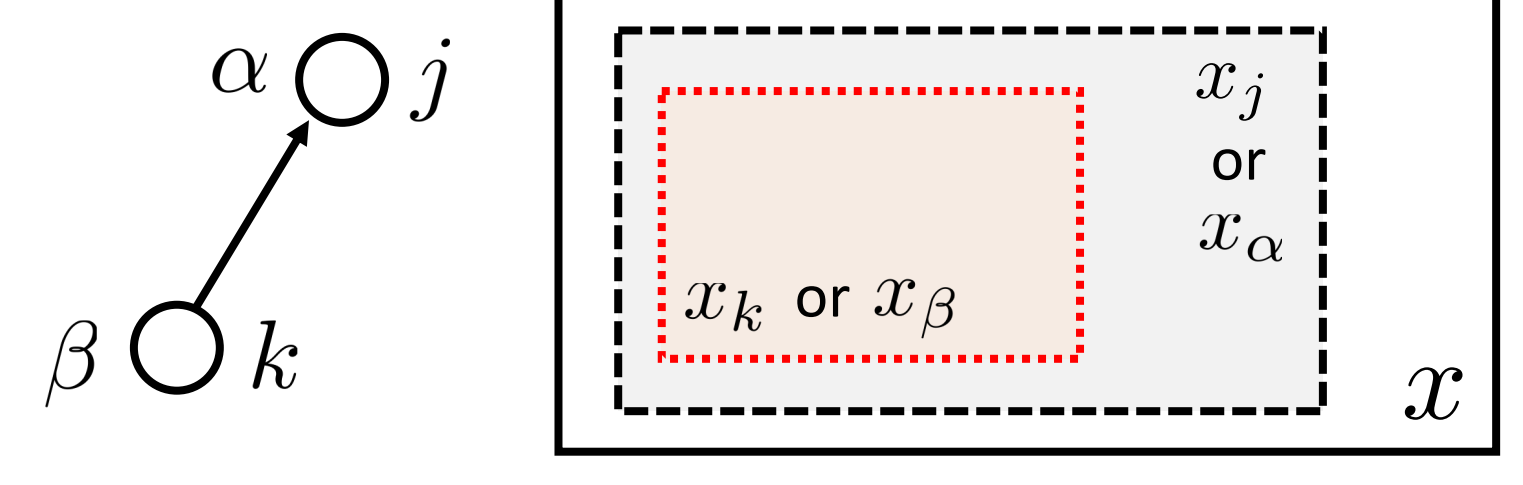
$f_j(z_\alpha) = f'_j(z_\alpha) \sum_{k \in \text{ch}(j)} w_{jk} \mathbb{E}_{z_\beta|z_\alpha} [f_k(z_\beta)]$

$g_k(z_\beta) = f'_k(z_\beta) \sum_{j \in \text{pa}(k)} w_{jk} \mathbb{E}_{z_\alpha|z_\beta} [g_j(z_\alpha)]$

Assumption

$\mathbb{P}(x_j|z_\alpha, z_\beta) = \mathbb{P}(x_j|z_\alpha)$

$\mathbb{P}(x_k|z_\alpha, z_\beta) = \mathbb{P}(x_k|z_\beta)$

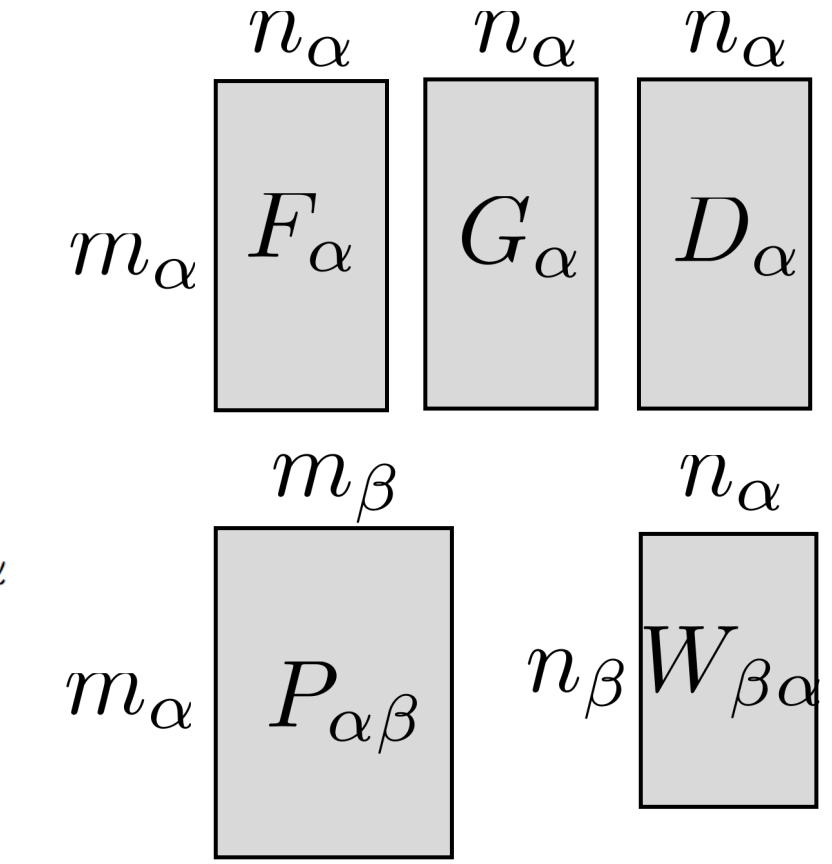


Matrix Form

$F_\alpha = D_\alpha \circ \sum_{\beta \in \text{ch}(\alpha)} P_{\alpha\beta} F_\beta W_{\beta\alpha}$

$\tilde{G}_\beta = D_\beta \circ \sum_{\alpha \in \text{pa}(\beta)} P_{\alpha\beta}^T \tilde{G}_\alpha W_{\beta\alpha}^T$

$\Delta W_{\beta\alpha} = (P_{\alpha\beta} F_\beta)^T \tilde{G}_\alpha$



n_α Number of nodes in the receptive field

m_α Number of event

$m_\alpha \rightarrow +\infty$

It becomes sample dimension

Geometrical Interpretation of Batch Norm

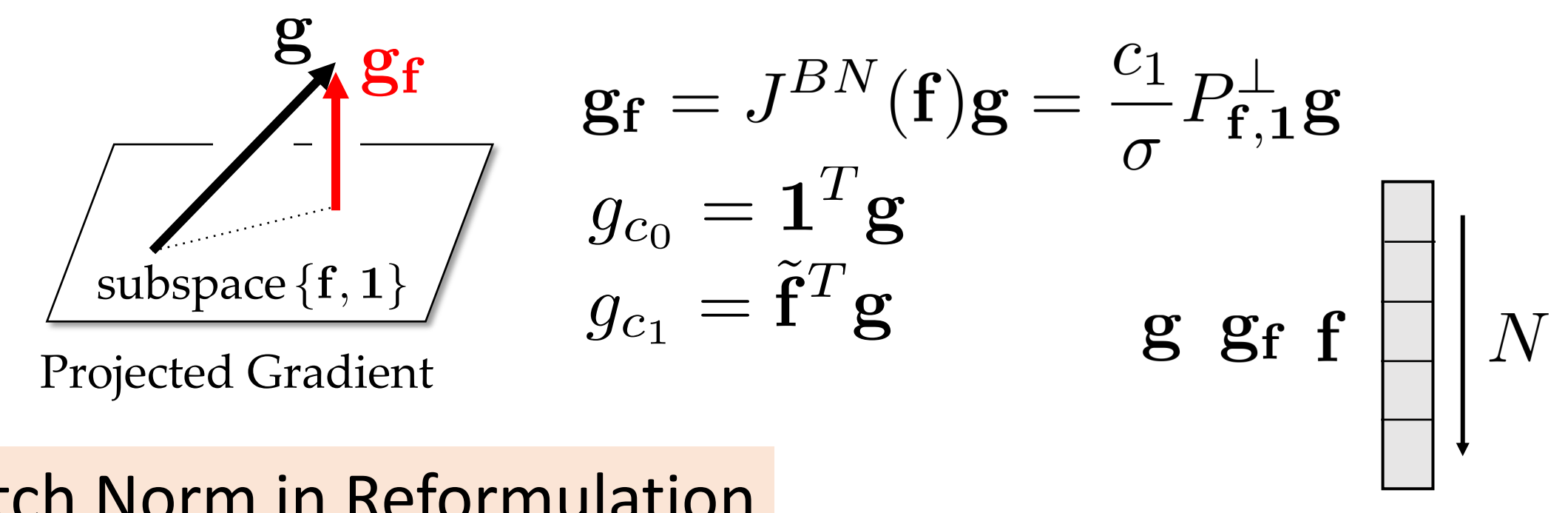
Forward Graph of Batch Norm

$\hat{\mathbf{f}} = P_1^\perp \mathbf{f}$, $\tilde{\mathbf{f}} = \hat{\mathbf{f}} / \|\hat{\mathbf{f}}\|_{\text{uni}}$, $\bar{\mathbf{f}} = c_1 \tilde{\mathbf{f}} + c_0 \mathbf{1} = S(\mathbf{f})\mathbf{c}$

$BN(\mathbf{f}) = c_1 \frac{P_1^\perp \mathbf{f}}{\|P_1^\perp \mathbf{f}\|_{\text{uni}}} + c_0 \mathbf{1}$

$\langle f, g \rangle_{\text{uni}} = \frac{1}{N} \sum_{i=1}^N f_i g_i$

Backward Gradient of Batch Norm [Theorem 5]



Backward Gradient from the original paper

$\frac{\partial \ell}{\partial \hat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \gamma$

$\frac{\partial \ell}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot (x_i - \mu_B) \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \cdot \frac{-1}{2(\sigma_B^2 + \epsilon)^{3/2}}$

$\frac{\partial \ell}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial \ell}{\partial \sigma_B^2} \cdot \frac{\sum_{i=1}^m -2(x_i - \mu_B)}{m}$

$\frac{\partial \ell}{\partial x_i} = \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial \ell}{\partial \sigma_B^2} \cdot \frac{2(x_i - \mu_B)}{m} + \frac{\partial \ell}{\partial \mu_B} \cdot \frac{1}{m}$

$\frac{\partial \ell}{\partial \gamma} = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i} \cdot \hat{x}_i$

$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i}$

Batch Norm in Reformulation

Using $\mathbb{E}_x [f_j(x)] = \mathbb{E}_{z_\alpha} [f_j(z_\alpha)]$ and we can rewrite BN:

$\mu = \mathbb{E}_{z_\alpha} [f_j]$, $\sigma^2 = \mathbb{E}_{z_\alpha} [(f_j(z_\alpha) - \mu)^2]$

$\mathbf{g}\mathbf{f} = J^{BN}(\mathbf{f})\mathbf{g} = \frac{c_1}{\sigma} P_{\mathbf{f}, \mathbf{1}}^\perp \mathbf{g}$

projection under new inner product $\langle f, g \rangle_{z_\alpha} = \mathbb{E}_{z_\alpha} [f(z_\alpha)g(z_\alpha)]$

Applications of the Framework

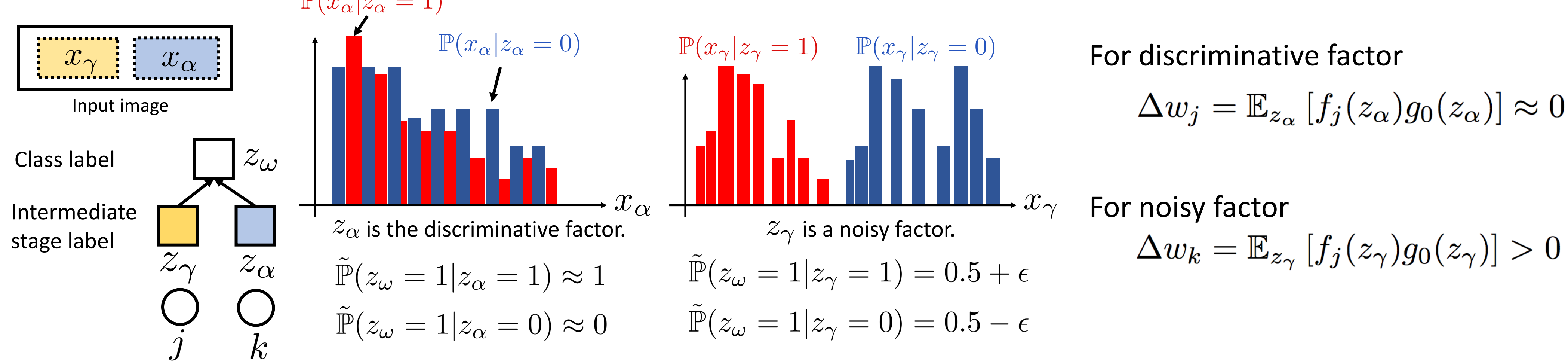
Linear versus Nonlinear

If $P_{\alpha\beta}$ is all-vert, $m_\alpha = n_\alpha = \mathcal{O}(\exp(\text{sz}(\alpha)))$, then [Theorem 6]

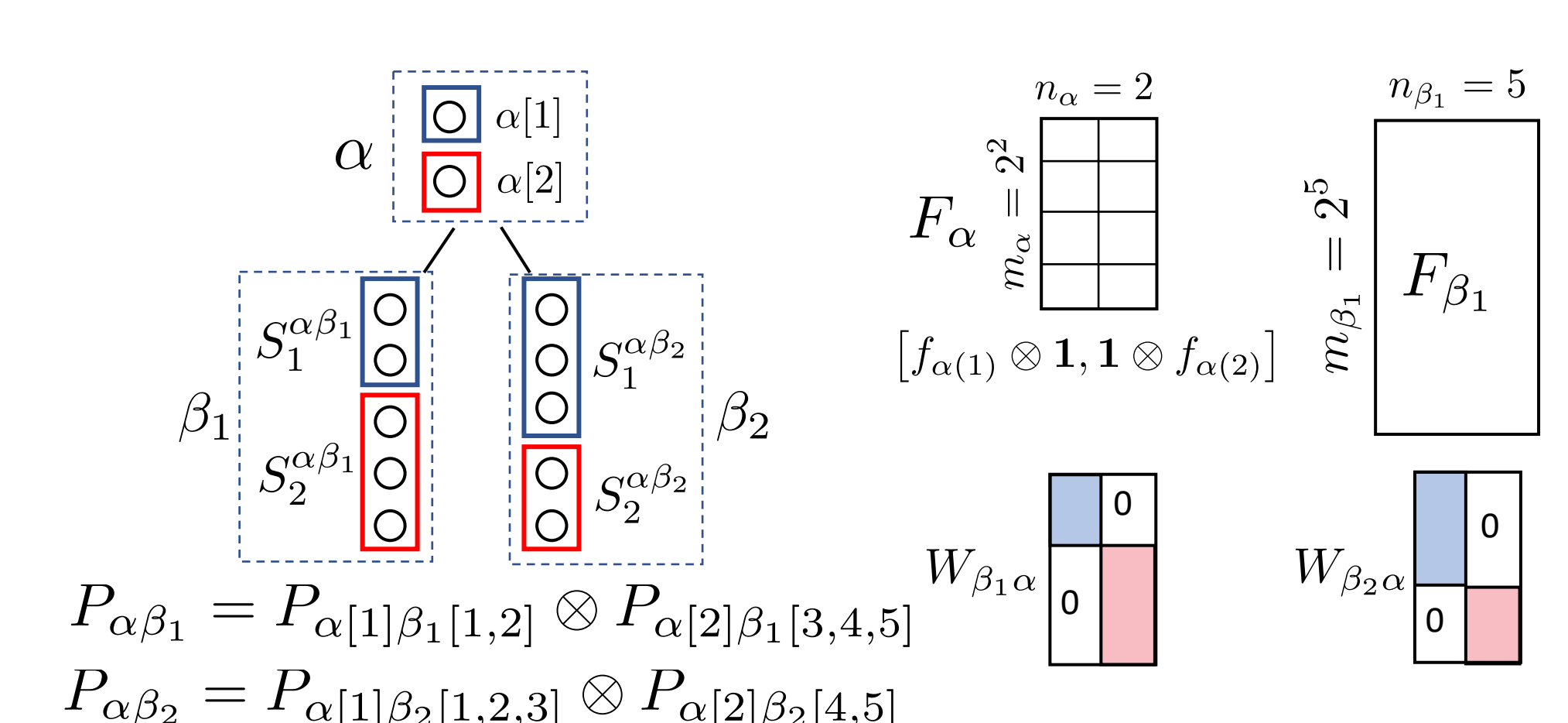
$\min_W \text{Loss}_{\text{ReLU}}(W) = 0$, $\min_W \text{Loss}_{\text{Linear}}(W) = \mathcal{O}(\exp(\text{sz}(\omega)))$

$\text{sz}(\alpha)$ Size of receptive field

Overfitting



Disentangled representation



Forward Disentanglement [Theorem 8]

If $P_{\alpha\beta}$ decomposes, F_α is disentangled, $W_{\beta\alpha}$ separates: Then F_β is also disentangled.

Weight update [Theorem 9]

If $P_{\alpha\beta}$ decomposes, F_β is disentangled, $W_{\beta\alpha}$ separates and G_α is also disentangled, then $\Delta W_{\beta\alpha}$ also separates.

Backward disentanglement

Still in progress ...