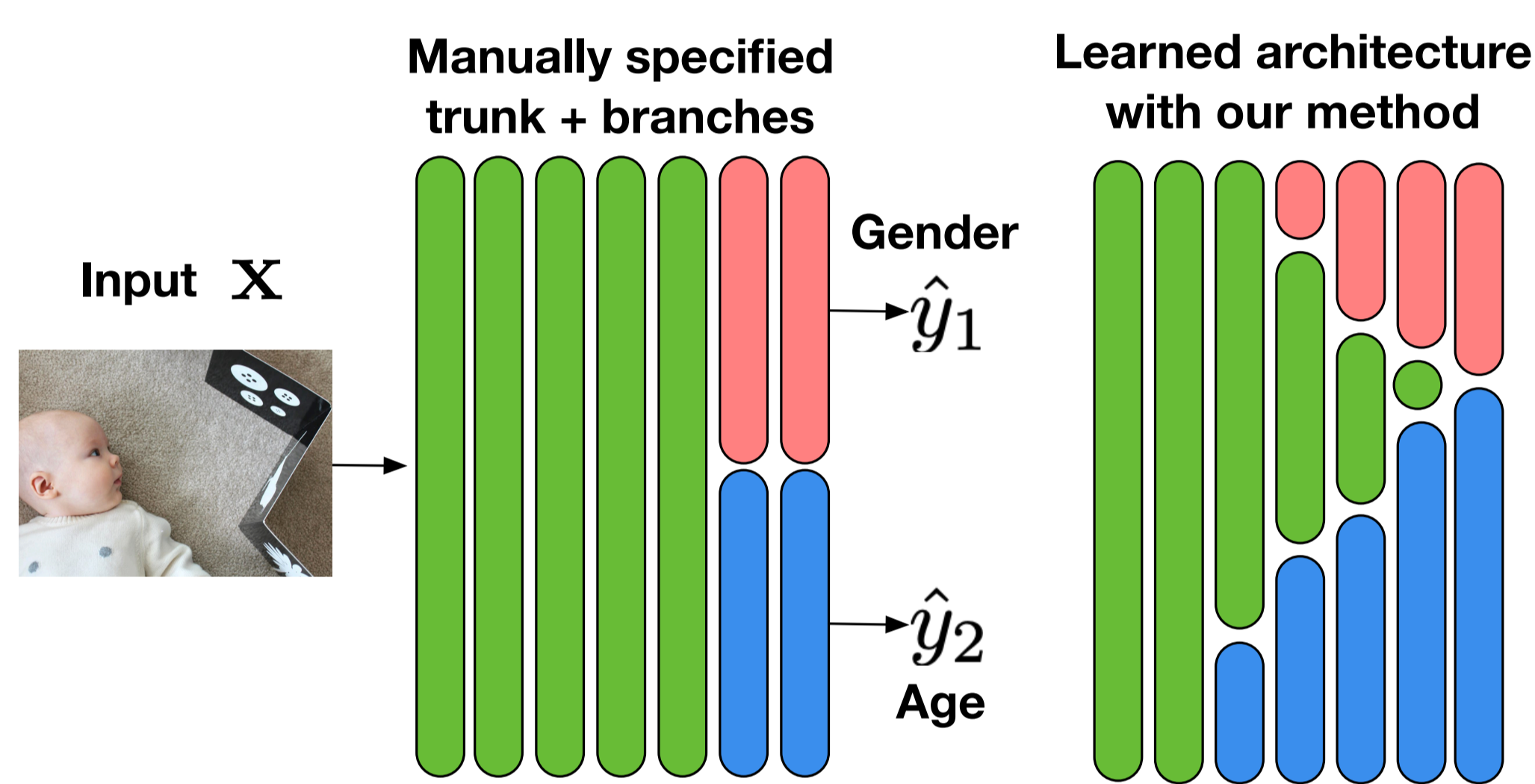




1 Summary

- Benefits of **multi-task learning (MTL)** depend on the nature of feature sharing and the network architecture
- Problem:** these architectures are **manually pre-specified** which can be **suboptimal**
- Solution:** we propose **Stochastic Filter Groups (SFG)**; a principled mechanism to learn the amount of feature sharing and separation between tasks
- We show the benefits of SFGs in two multi-task problems



2 What are Stochastic Filter Groups?

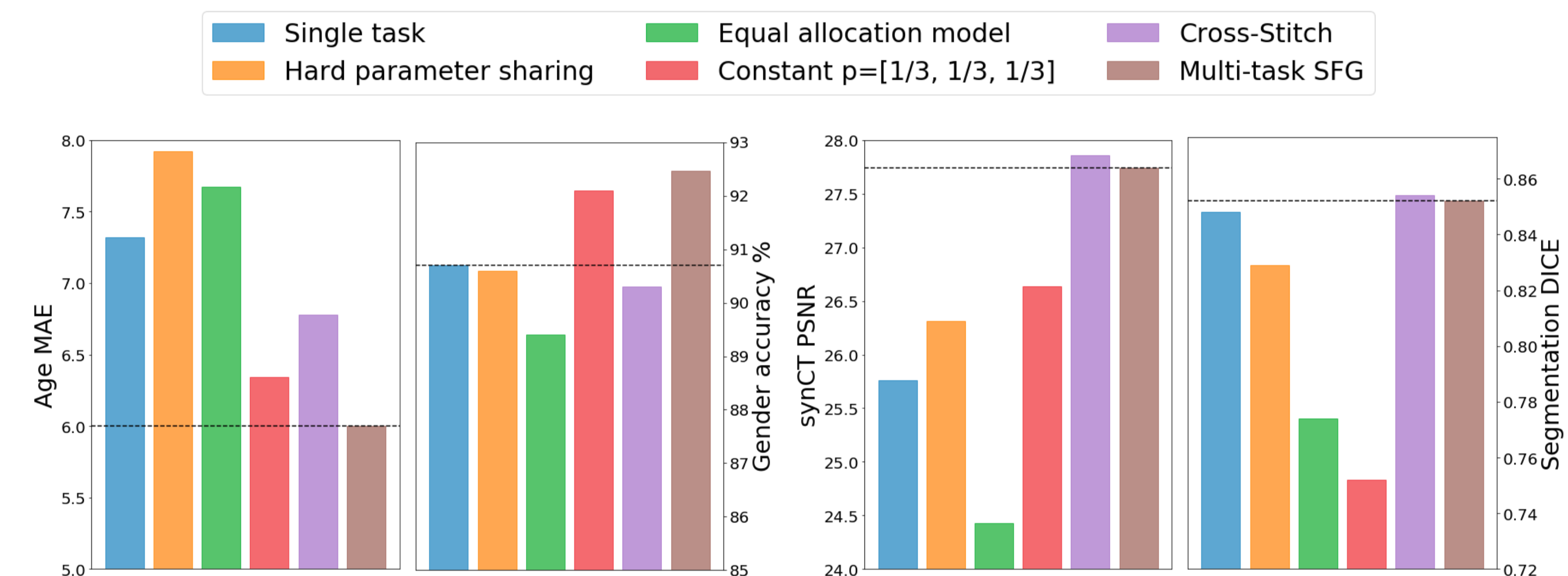
- Core idea:** cluster convolution kernels into task specific and shared groups in each layer of a CNN
- We define **task-specific groups** as the set of filters that are only updated to minimise corresponding task losses, while the **shared group** follows the same logic but is learned to optimise **all** task
- Jointly learn the grouping and convolution kernel weights

Filters	Group probabilities	Sample & Assign to Groups		
		p_1	p_s	p_2
w_1	Cat (0.6 0.1 0.3)	~	$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$	G_1 "Task 1"
w_2	Cat (0.1 0.9 0)	~	$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$	G_s "Shared"
w_3	Cat (0.4 0.3 0.3)	~	$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$	G_1 "Task 1"
w_4	Cat (0.1 0.8 0.1)	~	$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$	G_s "Shared"
w_5	Cat (1.0 0 0)	~	$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$	G_1 "Task 1"
w_6	Cat (0.05 0.05 0.9)	~	$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$	G_2 "Task 2"

Definitions:
 G_1 := filters updated w.r.t task-1 loss
 G_s := filters updated w.r.t task-1 loss and task-2 loss
 G_2 := filters updated w.r.t task-2 loss

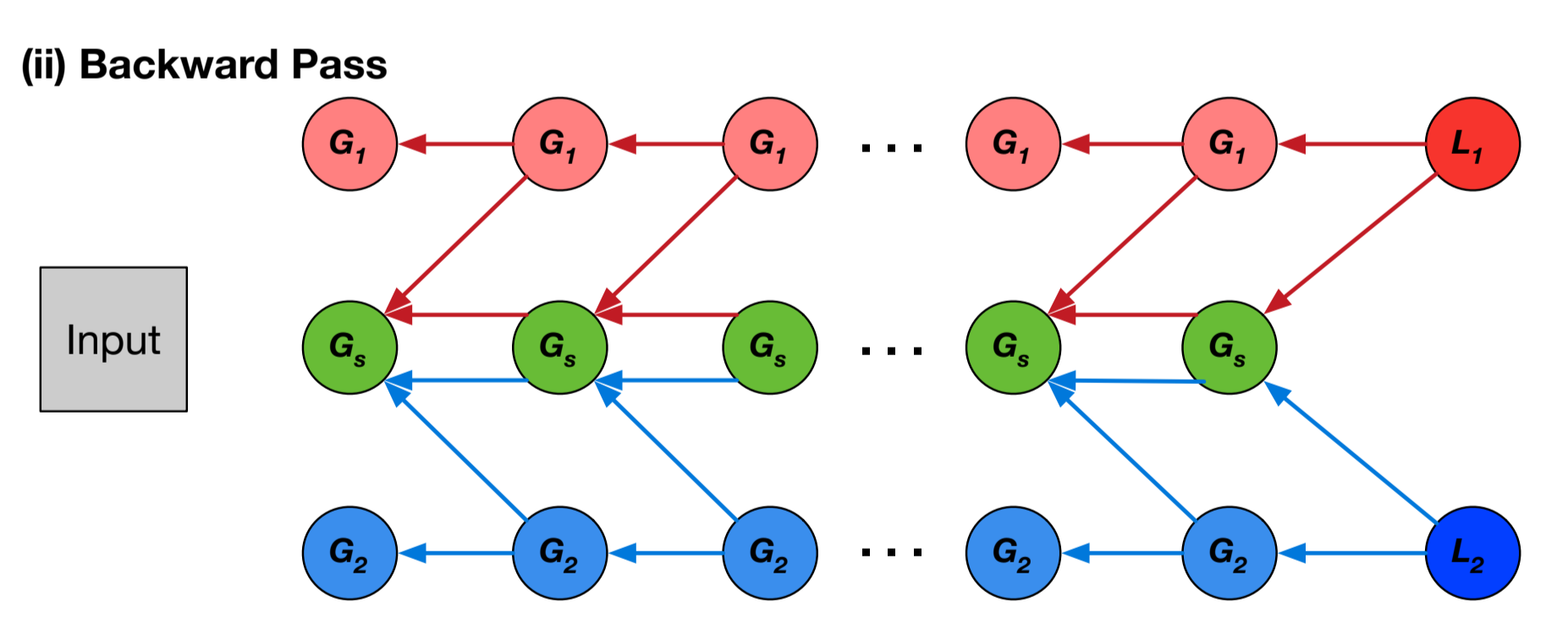
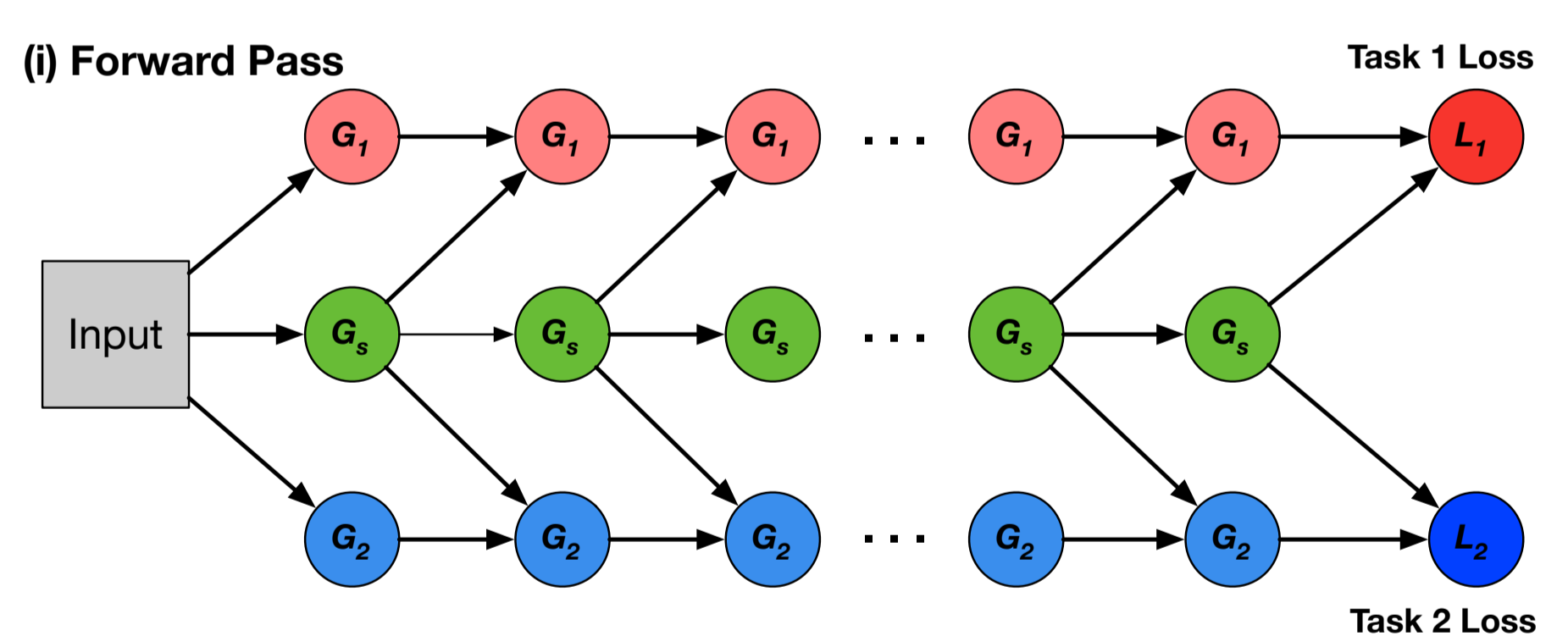
4 SFGs improve MTL performance

- Learning the allocation of kernels in MTL improves task performance
- We compared SFG-MTL architectures against: a) single-task networks, hard-parameter sharing networks, MTL networks with no learned allocation and Cross-Stitch networks [Misra et al, 2016]
- Dataset 1: UTKFace - age and gender prediction
- Dataset 2: Prostate MRI – CT synthesis and segmentation



3 Optimisation of network weights and kernel grouping

1 Sparse routing of features for desired gradient flow



2 Filter assignment as T+1 group drop-out

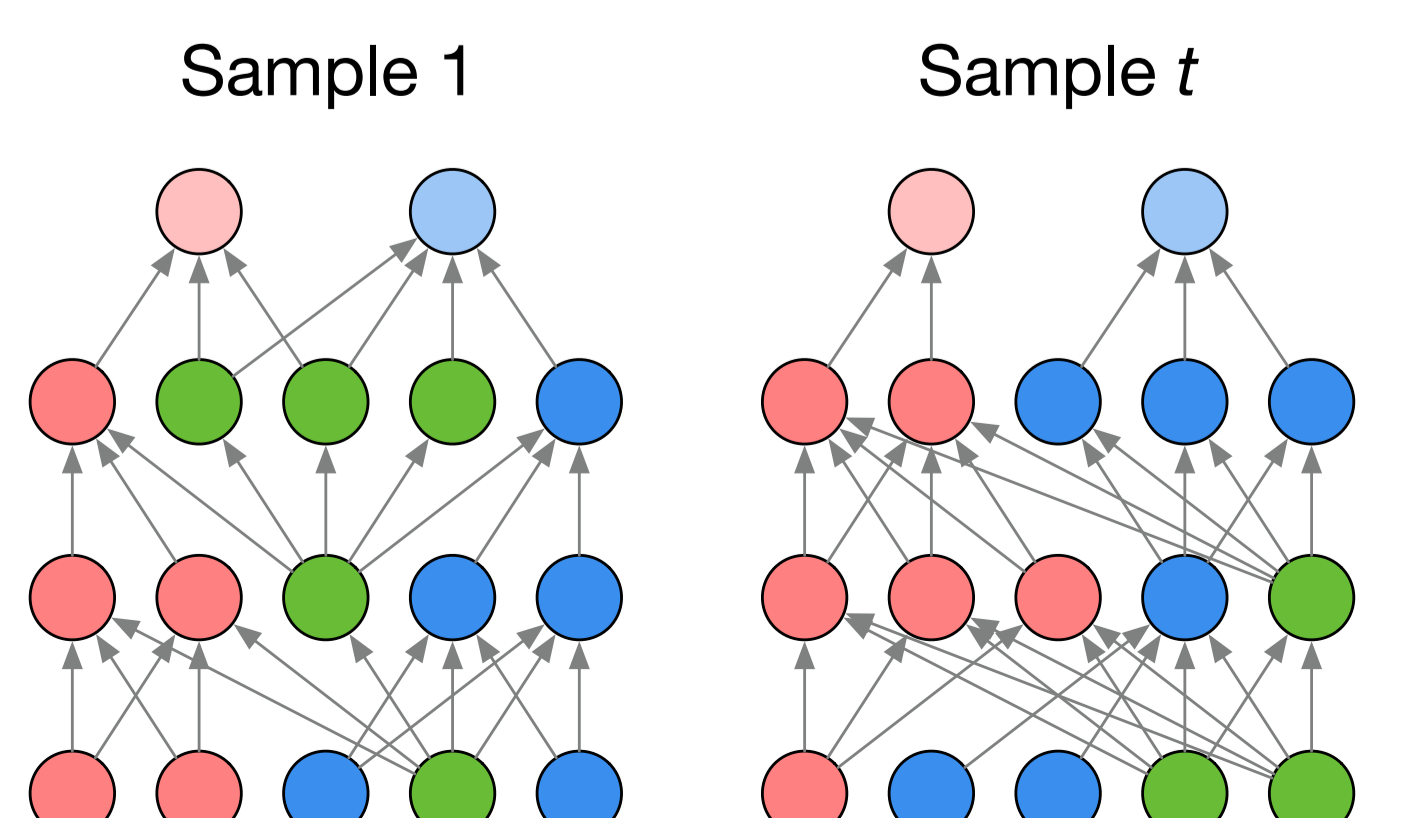
- Cast learning of grouping probabilities and filter weights as variational inference [Gal et al. 2018]
- Extended binary dropout to categorical distributions
- Minimise the following variational objective:

$$\mathcal{L}_{MC}(\phi) = -\frac{N}{M} \sum_{i=1}^M \left[\log p(y_i^{(1)} | \mathbf{x}_i, \mathcal{W}_i) + \log p(y_i^{(2)} | \mathbf{x}_i, \mathcal{W}_i) \right] + \lambda_1 \cdot \sum_{l=1}^L \sum_{k=1}^{K_l} \|\mathbf{M}^{(l,k)}\|^2 - \lambda_2 \cdot \sum_{l=1}^L \sum_{k=1}^{K_l} \mathcal{H}(\mathbf{p}^{(l,k)})$$

3 Continuous relaxation using Gumbel-Softmax [Jang et al. 2016]

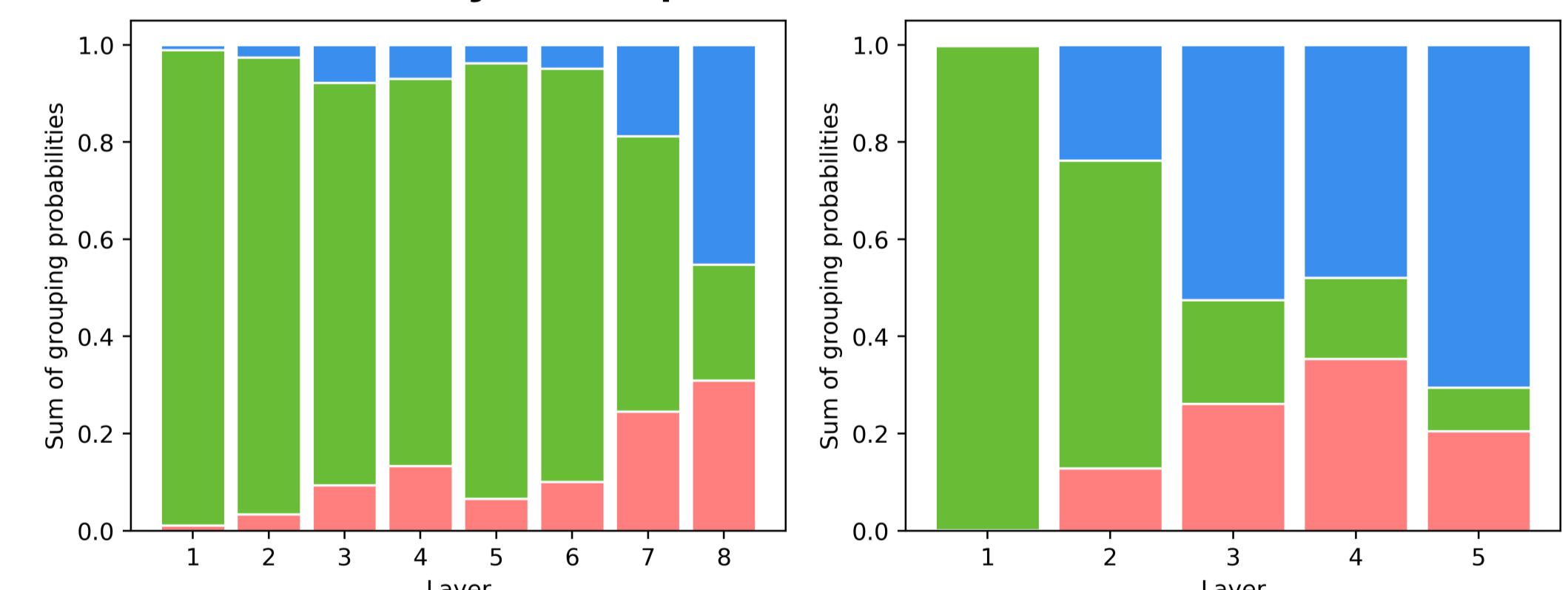
- First two terms have **zero** gradients w.r.t assignments probabilities \mathbf{p}

$$\text{Softmax}([g_i + \log p_i] / \tau) \quad g \sim \text{Gumbel}(0, 1)$$

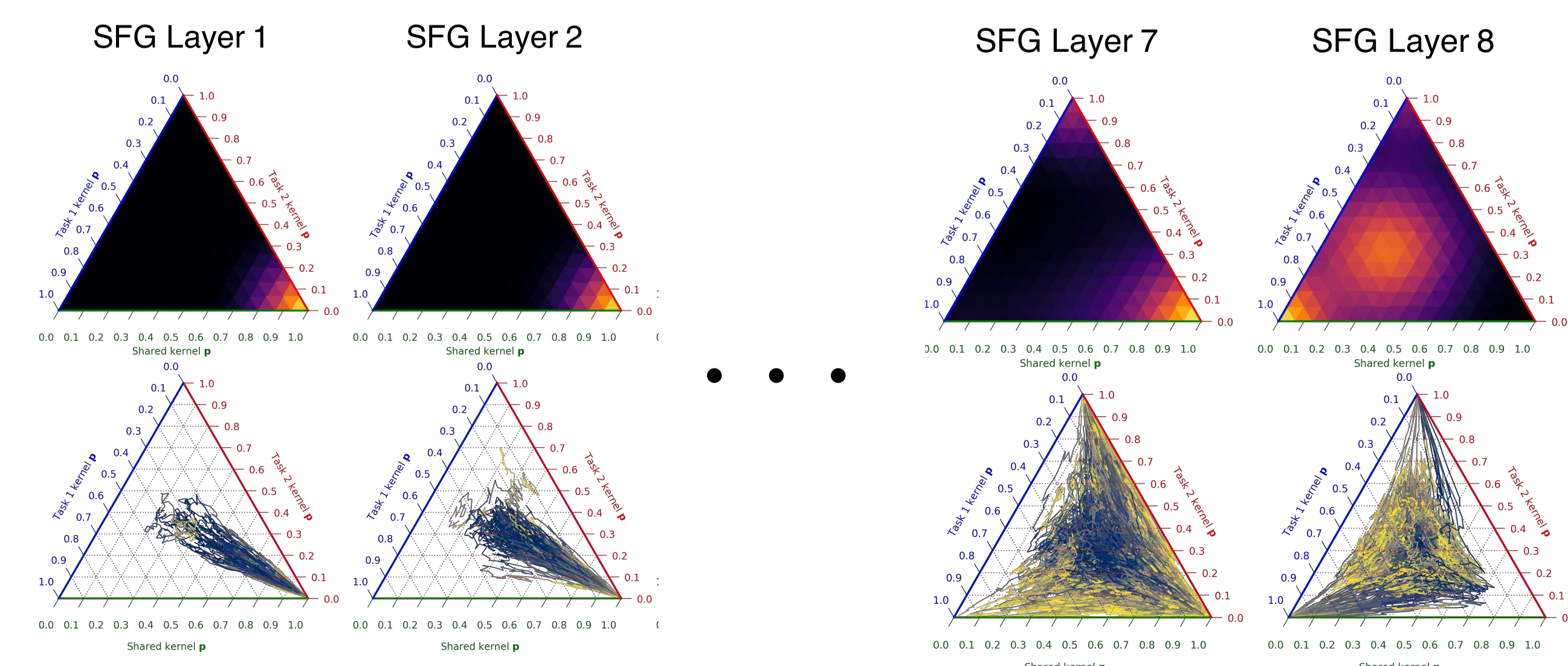


5 Qualitative Results

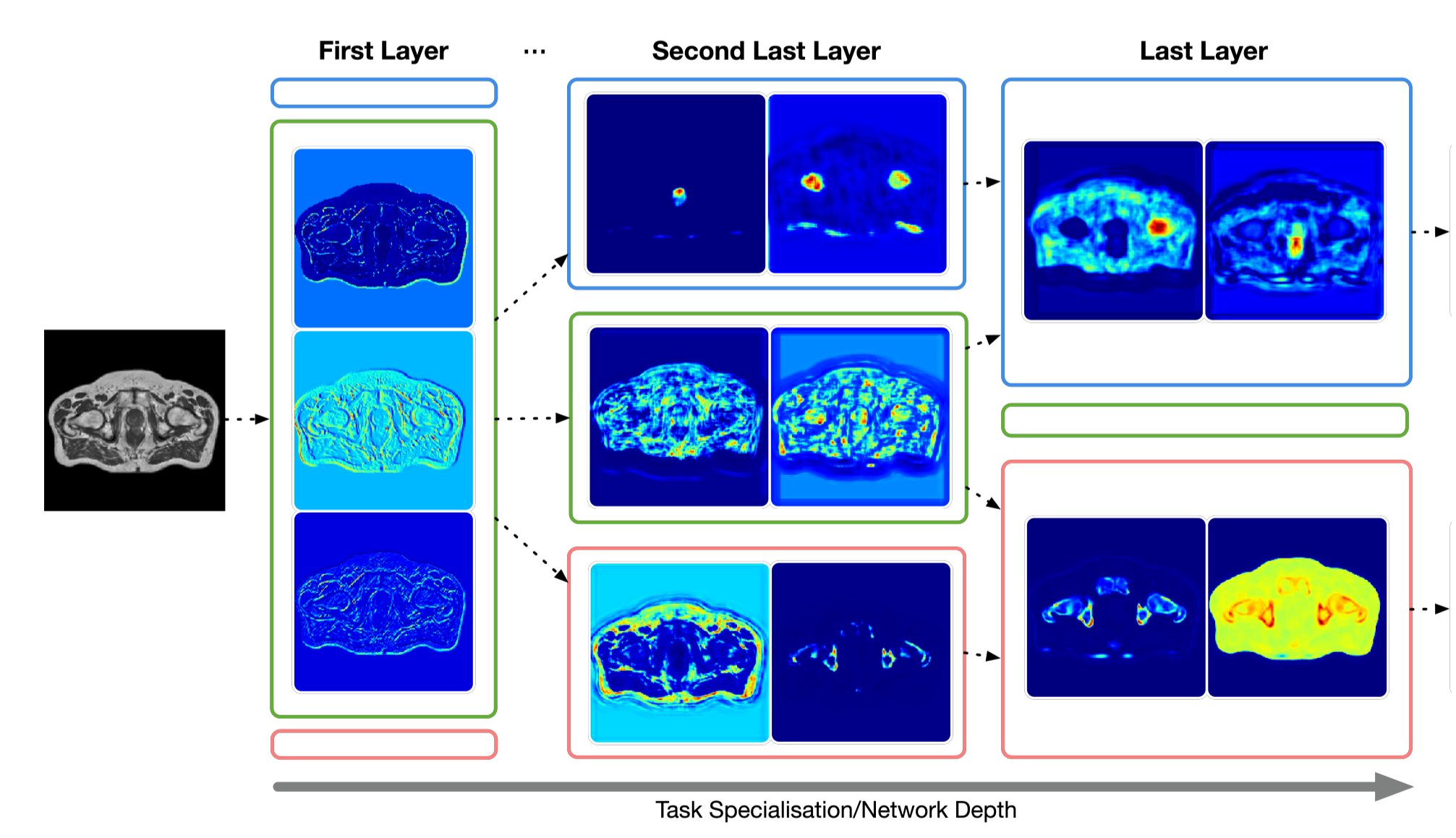
- We can visualise the learned allocation of kernels in a CNN with SFG modules
- Across both datasets, the ratio of task-specific groups increases with layer depth



- Density plots of probabilities \mathbf{p} illustrate learned grouping
- Training trajectories reveal some kernels converge faster to corresponding groups



- Activation maps for kernels with low grouping entropy confirm increasing task specialisation



- Maps for kernels with high grouping entropy show uncertainty in feature utility for maximising task performance

