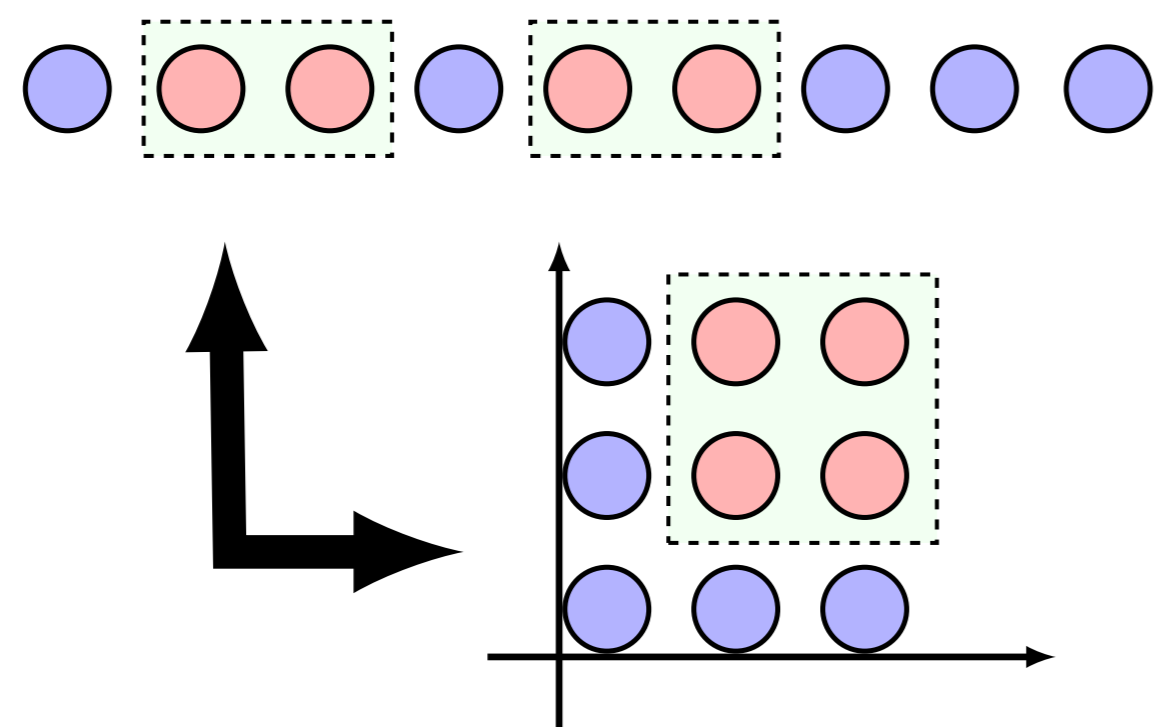


## 1. Introduction

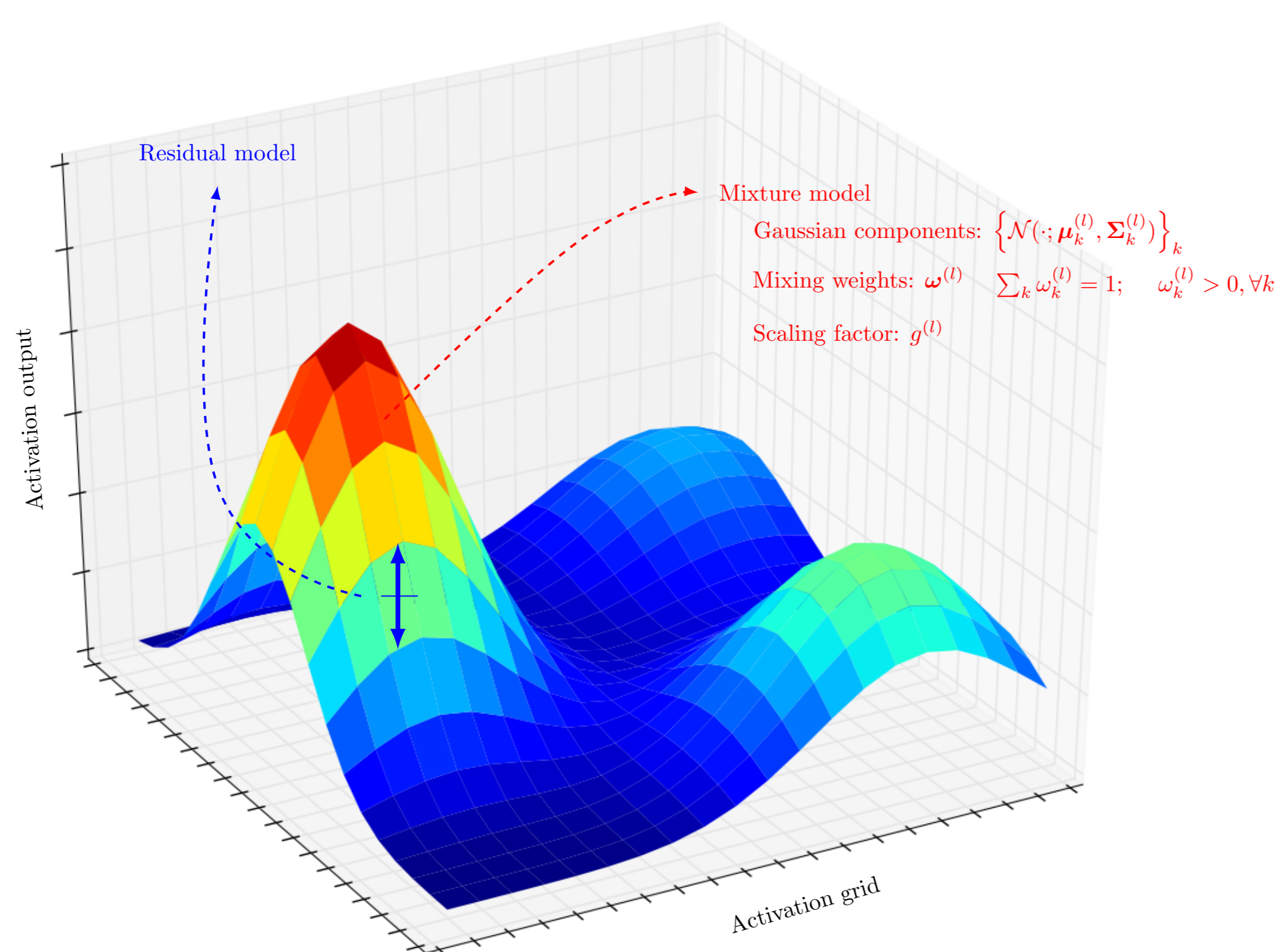
- ▶ DNNs are commonly treated as “black-box” models
  - ▶ Hard to interpret and group DNN parameters
  - ▶ Make regularisation and adaptation challenging
- ▶ Deep activation mixture models (DAMMs) are proposed to
  - ▶ Improve network regularisation and interpretation
  - ▶ Allow novel network adaptation schemes

## 2. Hidden Unit Reorganisation



- ▶ Reorganise units of each hidden layer to form a grid
  - ▶ Avoid the arbitrary ordering of hidden units
  - ▶ Enable activation functions to be related in regions of the network
- ▶ Each unit  $i$  is represented as a point  $s_i$  in the grid space

## 3. Network Topology



- ▶ How to relate different activation functions?

$$h^{(l)} = h_{\text{mix}}^{(l)} + h_{\text{res}}^{(l)}$$

- ▶ Mixture model — activation function contour

$$h_{\text{mix},i}^{(l)} = g^{(l)} \sum_{k=1}^K \omega_k^{(l)} \mathcal{N}(s_i; \mu_k^{(l)}, \Sigma_k^{(l)})$$

- ▶ Dynamic scaling factor  $g^{(l)}$  and mixing weights  $\omega^{(l)}$ 
  - $g^{(l)}$  — sigmoid       $\omega^{(l)}$  — softmax
- ▶ Interpret Gaussian component as phoneme or ...
  - ▶ E.g. set  $\{\mu_k^{(l)}\}$  as 2D projection of phoneme average features
  - ▶ Perform adaptation on Gaussian mean vector and covariance matrix

- ▶ Residual model — fluctuations on contour

$$h_{\text{res}}^{(l)} = \tanh(W^{(l)\top} h^{(l-1)} + c^{(l)})$$

- ▶ Add small “noise” to GMM contour

## 4. Training

- ▶ Highly restricted mixture model v.s. powerful residual model
- ▶ How to train mixture model to the maximal extent?

$$h^{(l)} \simeq h_{\text{mix}}^{(l)} \Rightarrow h_{\text{res}}^{(l)} \simeq \mathbf{0}$$

- ▶ Training criterion with residual-model regularisation

$$\mathcal{F}(\theta_{\text{mix}}, \theta_{\text{res}}) = \mathcal{L}(\theta_{\text{mix}}, \theta_{\text{res}}) + \eta \|\theta_{\text{res}}\|_2^2$$

- ▶ Isolating training mode

**Algorithm 1** Isolating Training Mode of DAMM.

```

1: for  $l := 1$  to  $L$  do
2:   initialise  $\theta_{\text{res}}^{(l)} = \mathbf{0}, \theta_{\text{mix}}^{(l)}$ 
3:   update  $\theta_{\text{mix}}$ 
4:   update  $\theta_{\text{res}}$ 
5: end for
6: finetune  $\theta_{\text{res}}$ 
    
```

## 5. Adaptation

- ▶ Adaptation on Gaussian components of mixture model

$$h_{\text{mix},i}^{(l)} = g^{(l)} \sum_{k=1}^K \omega_k^{(l)} \mathcal{N}(s_i; \mu_k^{(l)}, \Sigma_k^{(l)})$$

- ▶ A compact set of parameters to adapt  $\{\mu_k^{(l)}, \Sigma_k^{(l)}\}_{1 \leq l \leq L}$
- ▶ Allow robust and rapid adaptation

- ▶  $\Sigma_k^{(l)}$  can be rewritten using unit variance and correlation coefficient

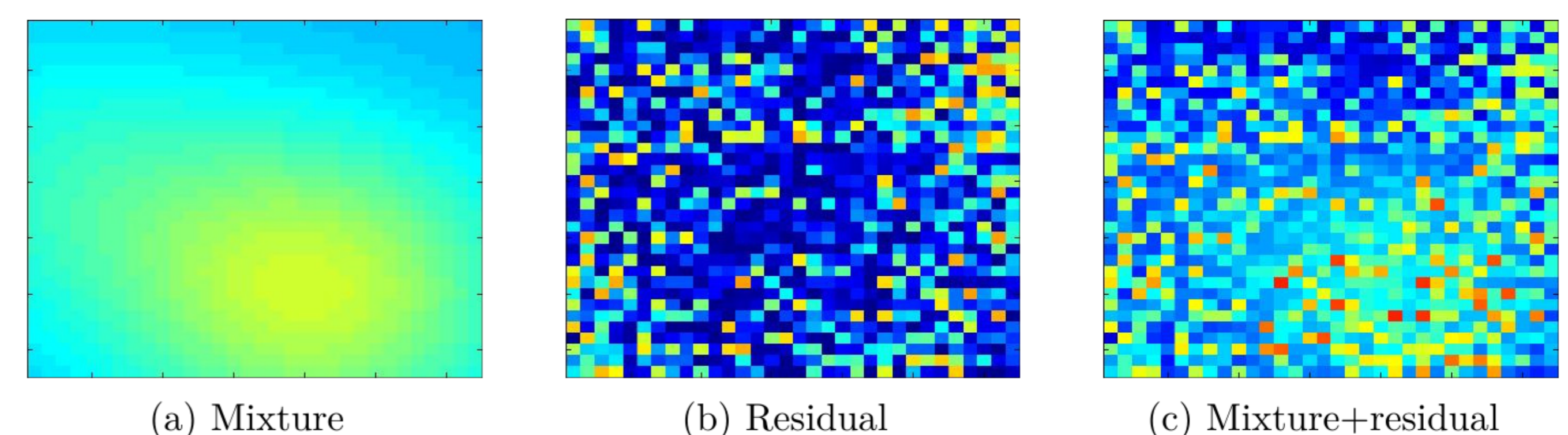
$$\Sigma_k^{(l)} = \begin{bmatrix} \sigma_{k1}^{(l)2} & \rho_k^{(l)} \sigma_{k1}^{(l)} \sigma_{k2}^{(l)} \\ \rho_k^{(l)} \sigma_{k1}^{(l)} \sigma_{k2}^{(l)} & \sigma_{k2}^{(l)2} \end{bmatrix}$$

- ▶ Parametrise  $\sigma_k^{(l)}$  and  $\rho_k^{(l)}$  to satisfy the positive-definite constraint

$$\sigma_k^{(l)} = \exp(\tilde{\sigma}_k^{(l)}), \rho_k^{(l)} = \tanh(\tilde{\rho}_k^{(l)})$$

## 6. Experiment

- ▶ Data and setup
  - ▶ 144-hour English broadcast news dataset (LDC97S44, LDC98S71)
  - ▶ DNN-HMM hybrid ASR framework
  - ▶ 5 hidden layers with 1024 units for both DNN and DAMM systems
  - ▶ 46 Gaussian components to form DAMM mixture models
  - ▶ Unsupervised utterance-level adaptation
- ▶ Grid-output example of mixture and residual models



- ▶ CE performance comparison [WER (%)]

System	Dev	Eval
DNN(tanh)	12.8	11.0
DAMM	<b>12.3</b>	<b>10.6</b>
DNN(sigmoid)	12.4	10.8

- ▶ Adaptation of mean, variance and correlation coefficient on CE DAMM

System	Adapt			Dev	Eval
	mean	variance	correlation		
SI	✗	✗	✗	12.3	10.6
SD	✓	✗	✗	12.2	10.6
	✗	✓	✓	12.1	10.5
	✓	✓	✓	<b>12.0</b>	<b>10.4</b>

- ▶ More effective to adapt covariance matrix than mean vector

- ▶ MPE performance comparison

System	Adapt			Dev	Eval
	mean	variance	correlation		
DNN(sigmoid)	-			11.4	10.1
DAMM	✗	✗	✗	11.4	10.0
	✓	✓	✓	<b>11.1</b>	<b>9.8</b>

- ▶ Up to 3% rel. WER reduction by adaptation

## 7. Conclusions

- ▶ Propose DAMM for network regularisation and adaptation
- ▶ Extend L2 regularisation to approach a dynamic surface, not zero
- ▶ Novel adaptation scheme to modify the dynamic surface