

## Mathematical Model

- Multiple Kernel Learning (MKL) [1]

$$f(\mathbf{x}) = \sum_{m=1}^P \langle \mathbf{w}_m, \Phi_m(\mathbf{x}) \rangle + b$$

$$= \sum_{i=1}^N \alpha_i y_i \sum_{m=1}^P \eta_m \underbrace{\langle \Phi_m(\mathbf{x}_i), \Phi_m(\mathbf{x}) \rangle}_{k_m(\mathbf{x}_i, \mathbf{x})} + b$$

- Localized Multiple Kernel Learning (LMKL) [2]

$$f(\mathbf{x}) = \sum_{m=1}^P \eta_m(\mathbf{x}|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{x}) \rangle + b$$

$$= \sum_{i=1}^N \alpha_i y_i \underbrace{\sum_{m=1}^P \eta_m(\mathbf{x}_i|\mathbf{V}) k_m(\mathbf{x}_i, \mathbf{x}) \eta_m(\mathbf{x}|\mathbf{V})}_{k\eta(\mathbf{x}_i, \mathbf{x})} + b$$

## Primal Problem

$$\min. \frac{1}{2} \sum_{m=1}^P \|\mathbf{w}_m\|^2 + C \sum_{i=1}^N \xi_i$$

$$\text{w.r.t. } \mathbf{w}_m \in \mathbb{R}^{D_m}, b \in \mathbb{R}, \boldsymbol{\xi} \in \mathbb{R}_+^N, \mathbf{V} \in \mathbb{R}^{D_G}$$

$$\text{s.t. } y_i \left( \sum_{m=1}^P \eta_m(\mathbf{x}_i|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{x}_i) \rangle + b \right) \geq 1 - \xi_i \quad \forall i$$

- For a given  $\mathbf{V}$

## Dual Problem

$$\max. J(\mathbf{V}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k\eta(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

- $\frac{\partial J(\mathbf{V})}{\partial \mathbf{V}}$  is used to update  $\mathbf{V}$

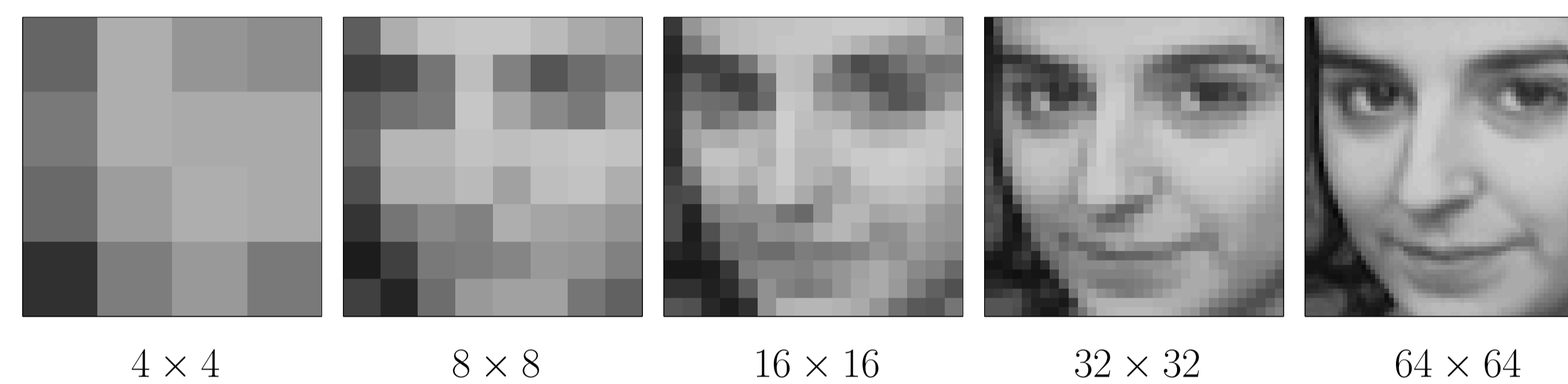
- Two-step alternate optimization algorithm for learning  $\boldsymbol{\alpha}$  and  $\mathbf{V}$

## Combining Multiple Resolutions

### Softmax Gating

$$\eta_m(\mathbf{x}|\mathbf{V}) = \frac{\exp(\langle \mathbf{v}_m, \Phi_G(\mathbf{x}) \rangle + v_{m0})}{\sum_{k=1}^P \exp(\langle \mathbf{v}_k, \Phi_G(\mathbf{x}) \rangle + v_{k0})} \quad \forall m$$

- Face image in different resolutions



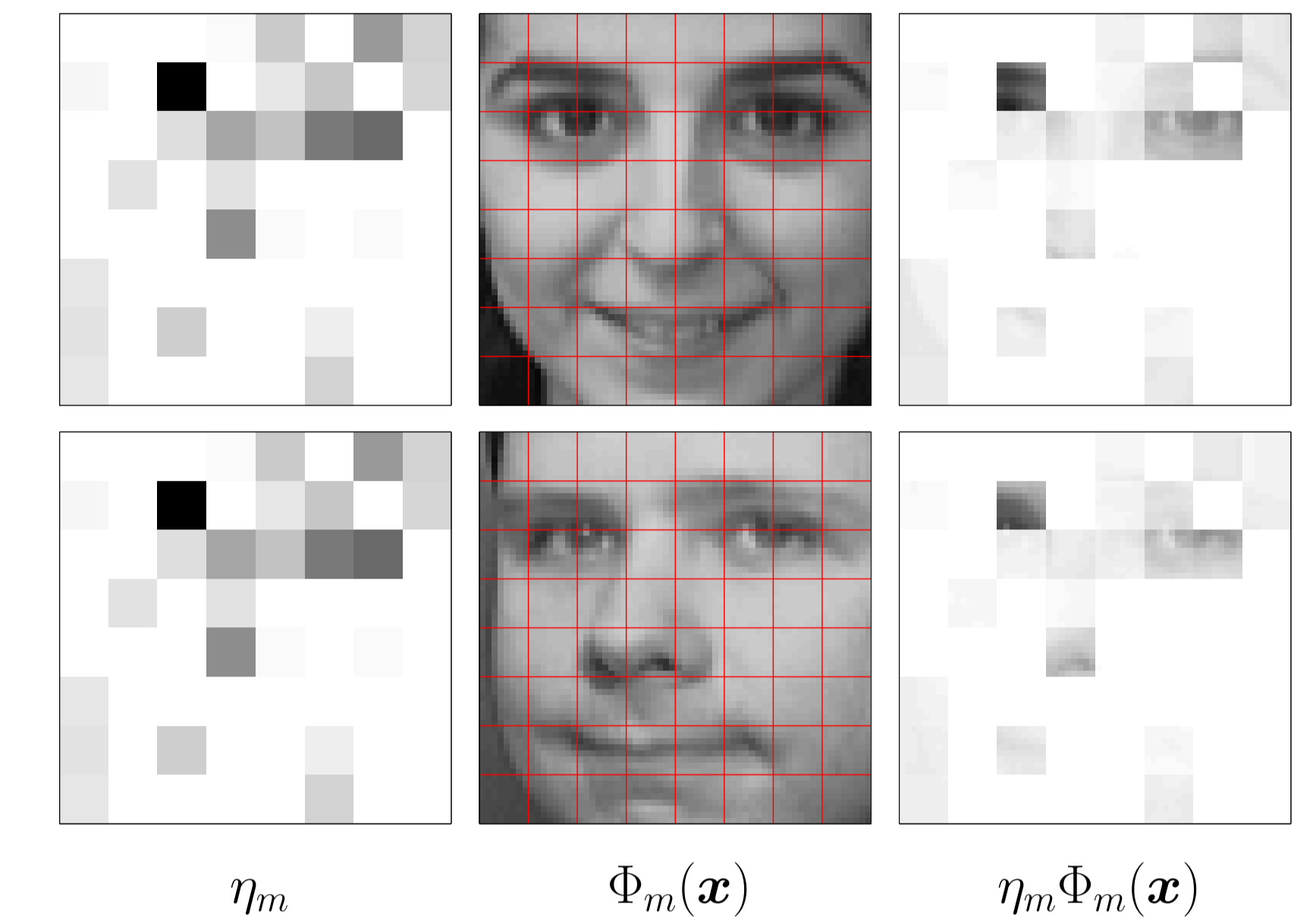
	Method	Configuration	Accuracy	SV
Single Kernel	SVM	$\Phi(\mathbf{x}) = 4 \times 4$	93.28 ± 0.65	21.70 ± 0.93
	SVM	$\Phi(\mathbf{x}) = 8 \times 8$	97.50 ± 1.16	20.13 ± 1.04
	SVM	$\Phi(\mathbf{x}) = 16 \times 16$	97.03 ± 0.93	19.82 ± 0.94
	SVM	$\Phi(\mathbf{x}) = 32 \times 32$	97.97 ± 1.48	23.71 ± 1.39
	SVM	$\Phi(\mathbf{x}) = 64 \times 64$	97.66 ± 1.41	25.94 ± 1.01
Multiple Kernel	MKL	uses only 64 × 64	97.66 ± 1.41	25.94 ± 1.01
	LMKL	$\Phi_G(\mathbf{x}) = 4 \times 4$	97.03 ± 1.15	29.29 ± 2.90
	LMKL	$\Phi_G(\mathbf{x}) = 8 \times 8$	99.38 ± 0.94	27.68 ± 2.95
	LMKL	$\Phi_G(\mathbf{x}) = 16 \times 16$	98.59 ± 1.41	26.52 ± 2.37
	LMKL	$\Phi_G(\mathbf{x}) = 32 \times 32$	99.38 ± 1.16	24.78 ± 2.57
LMKL	$\Phi_G(\mathbf{x}) = 64 \times 64$	99.53 ± 0.65	26.65 ± 4.12	

## Combining Multiple Input Patches

### Sigmoid Gating

$$\eta_m(\mathbf{x}|\mathbf{V}) = \frac{1}{1 + \exp(-\langle \mathbf{v}_m, \Phi_G(\mathbf{x}) \rangle - v_{m0})} \quad \forall m$$

Method	Configuration	Accuracy	SV
MKL	$\Phi_m(\mathbf{x}) = 8 \times 8$	99.38 ± 0.94	19.42 ± 0.87
LMKL	$\Phi_m(\mathbf{x}) = 8 \times 8 \quad \Phi_G(\mathbf{x}) = 8 \times 8$	99.84 ± 0.44	24.38 ± 1.96
MKL	$\Phi_m(\mathbf{x}) = 16 \times 16$	99.06 ± 0.88	22.19 ± 1.00
LMKL	$\Phi_m(\mathbf{x}) = 16 \times 16 \quad \Phi_G(\mathbf{x}) = 4 \times 4$	99.53 ± 0.65	23.35 ± 1.47



## Conclusions

- Gating model acts as a saliency detector for selective attention
- Use a cheap gating model that selects among kernels of different costs only when needed

## References

- [1] F. R. Bach, G. R. G. Lanckriet, and M. I. Jordan. Multiple kernel learning, conic duality, and the SMO algorithm. In *Proceedings of the 21st International Conference on Machine Learning*, pages 41–48, 2004.
- [2] M. Gönen and E. Alpaydın. Localized multiple kernel learning. In *Proceedings of the 25th International Conference on Machine Learning*, pages 352–359, 2008.