# Learning with Whom to Share in Multi-task Feature Learning

Zhuoliang Kang

Department of Computer Science

University of Southern California

[Joint work with Kristen Grauman @ U. of Texas Austin and Fei Sha @ U. of Southern California]

# This is a talk about

- ***Multi-task learning.***
- ***Automatic tasks grouping.***

E.g. multiple animal recognition tasks.



persian cat

horse
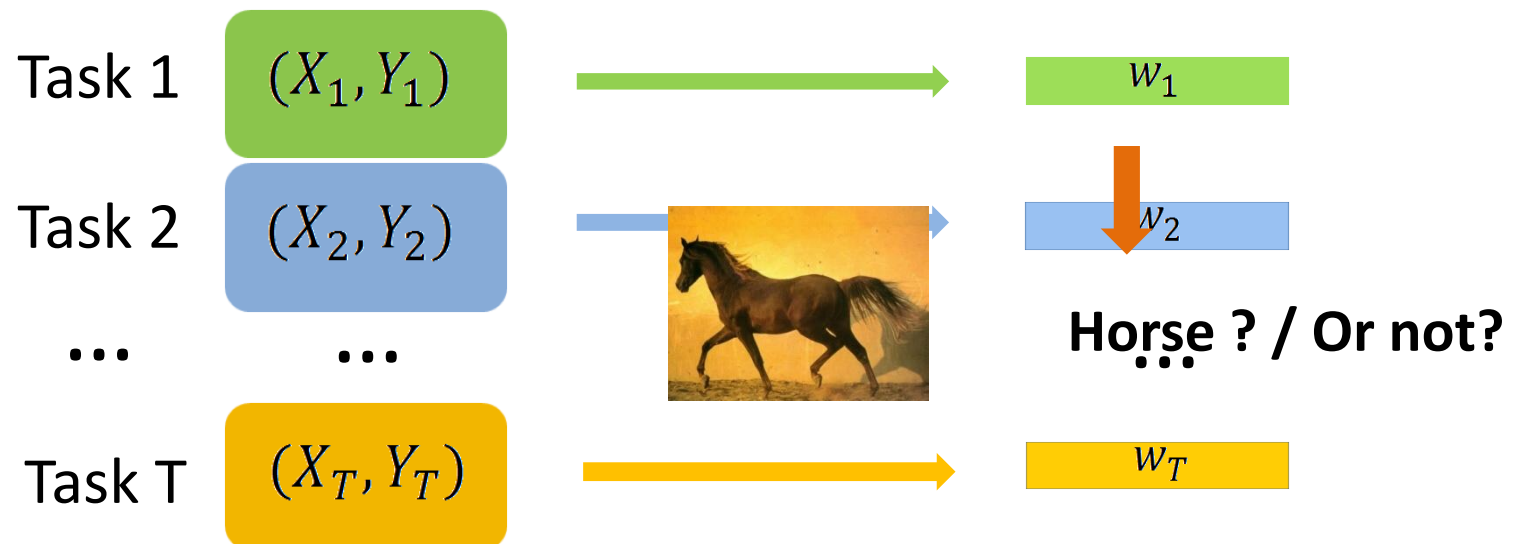
bobcat

buffalo

ox

siaseme cat

# Outline

- **Background**
  - What is multi-task learning
- **Motivation**
  - Why we want to group tasks
- **Algorithms**
  - How to discover the grouping
- **Empirical results**
  - Validate our approach
- **Conclusion**
  - Summary
  - Future work

# Supervised learning

- Given training data and label
  - Learn parameters for future prediction.
- Given **multiple** tasks.
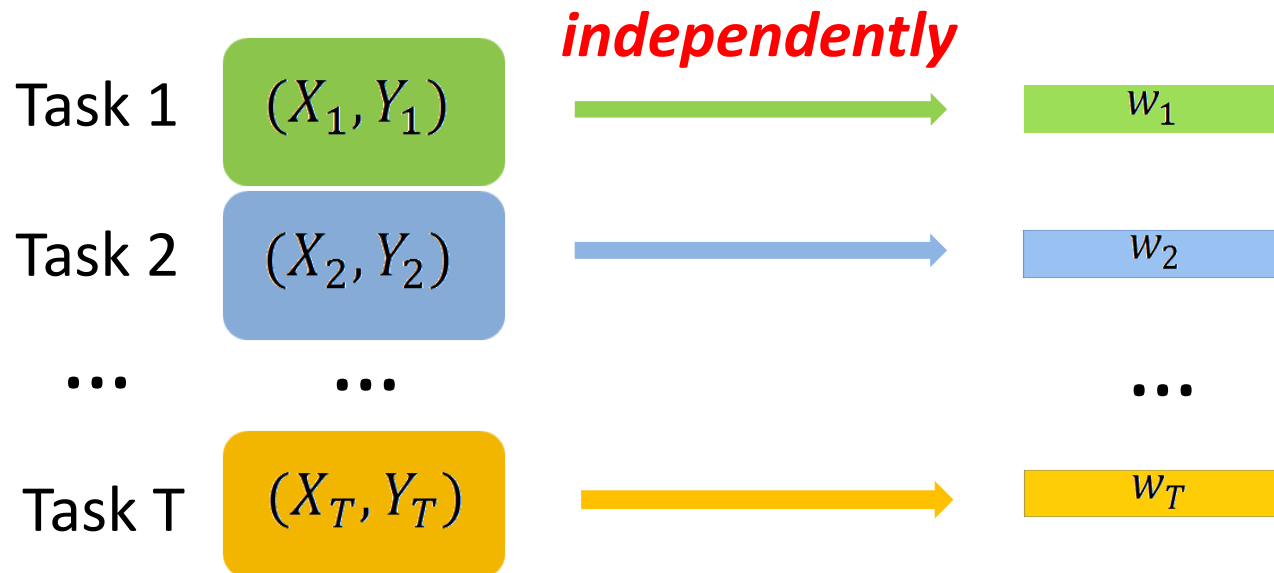  - Learn parameters **independently**.



Task 1 $(X_1, Y_1)$ $\longrightarrow$ $w_1$

Task 2 $(X_2, Y_2)$ $w_2$

... ... **Horse ? / Or not?**

...

Task T $(X_T, Y_T)$ $\longrightarrow$ $w_T$

# Regularization based framework

For ***each task***, solve an optimization problem

Balance empirical risk and model complexity

$$\min_{w_t} \quad loss(w_t, X_t, Y_t) + R(w_t)$$

***independently***

Task 1 $(X_1, Y_1)$ $\longrightarrow$ $w_1$

Task 2 $(X_2, Y_2)$ $\longrightarrow$ $w_2$

... ... ...

Task T $(X_T, Y_T)$ $\longrightarrow$ $w_T$

# How to solve a group of related tasks?

- Example
  - Recognizing similar animals.
  - Recognizing similar handwritten digits.
- We can do better than learning independently.



*independently*

$(X_1, Y_1) \longrightarrow w_1$

$(X_2, Y_2) \longrightarrow w_2$

$\cdots$

$(X_T, Y_T) \longrightarrow w_T$

# Multi-task learning (MTL)

- Main idea
  - Learn multiple tasks *jointly*.
  - Take the advantage of *relatedness*.
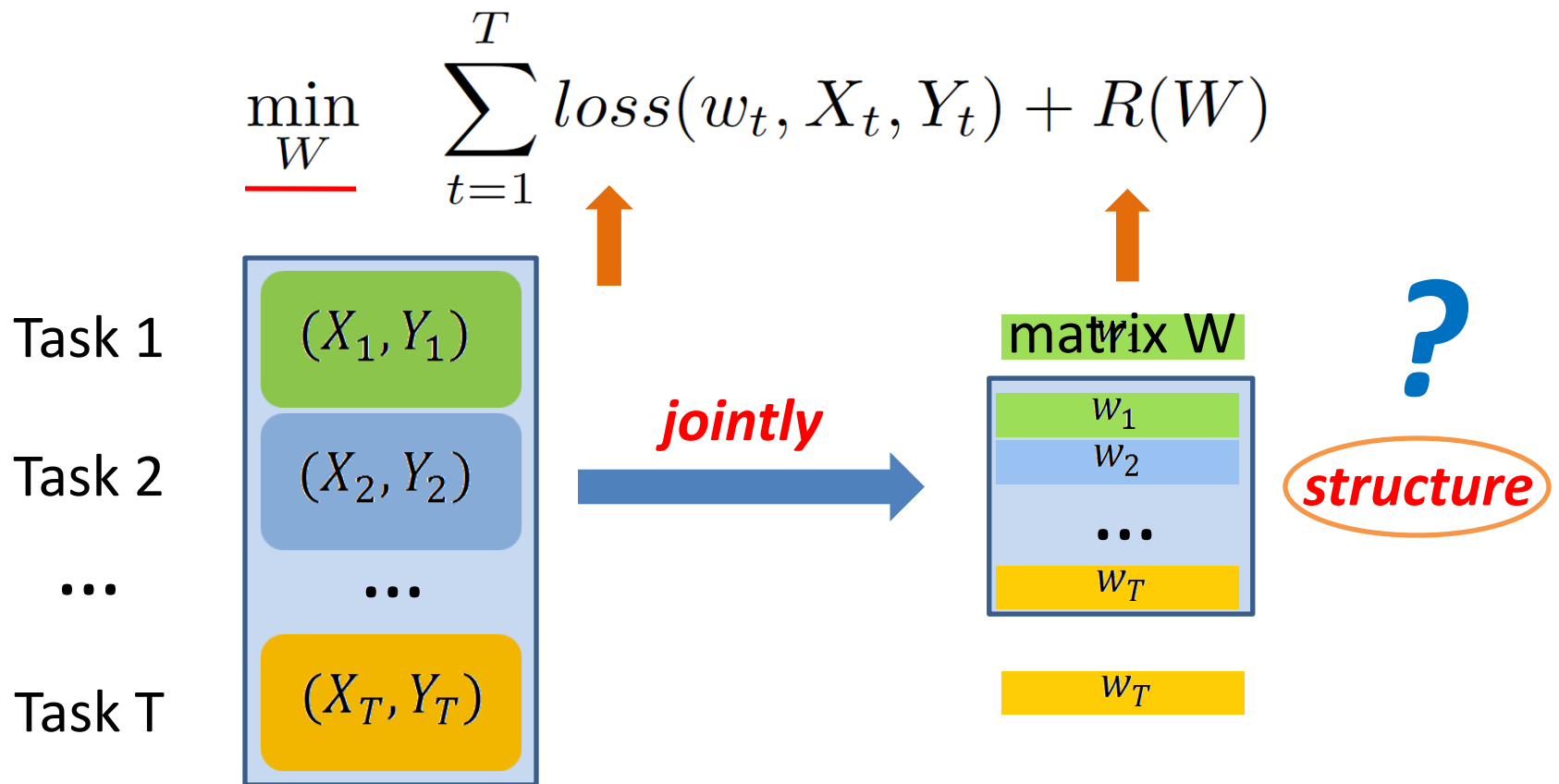
- Benefits
  - Improve *generalization* performance.
  - Require *less* training data.

**Related work**

[ Caruana, 97. Bakker and Heskes, 03. Evgeniou, et al. 04. Ando and Zhang. 05. Yu, et al., 05. Lee, et al., 07. Argyriou, et al. 08, Daumé, 09. Parameswaran, S. and Weinberger, K.Q. 2009. … ]

# Regularization based approach

Solve a joint optimization problem *for all tasks*.

Balance between *total empirical risk* and *relatedness*.

$$\min_{W} \sum_{t=1}^{T} loss(w_t, X_t, Y_t) + R(W)$$

Task 1 $(X_1, Y_1)$

Task 2 $(X_2, Y_2)$

$\cdots$ $\cdots$

Task T $(X_T, Y_T)$

*jointly*

matrix W

$w_1$
$w_2$
$\cdots$
$w_T$

**?**

*structure*

$w_T$

[Evgeniou, et al. 2004. Parameswaran, S. and Weinberger, K.Q. 2009. Zhang, et al. 2010. ...]

# Alternatives to regularization based MTL

- Share a common layer in Neural Network
  - R. Caruana, 1997.
  - B. Bakker and T. Heskes, 2003.
- Share common priors
  - Yu, et al., 2005.
  - Lee, et al., 2007.
  - E. Bonilla, et al. 2008
  - Daumé, III, Hal. 2009.
- etc …

# Multi-task feature learning (MTFL)

[ Argyriou, et al. 2008. ]

Task-relatedness

- Parameters lie on *a common low-dimensional subspace.*
- Or equivalently, models share *a common feature subspace.*

matrix W                    subspace

*structure*



Structural constraint on W: **low rank**

# **Low-rank Regularization**

sub-space

$w_1$    $w_2$

$w_T$

*Rank*: number of none-zero singular values    **( non-convex )**

$$W = U \begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_D \end{pmatrix} V^T \qquad \mathrm{rank}(W) = \sum_d 1[\sigma_d \neq 0]$$

Singular Value Decomposition

Convex relaxation

– *Trace norm*: $L_1$-norm of singular values    **( convex )**

$$\|W\|_{tr} = \sum_d |\sigma_d|$$

# Outline

- **Background**
  - What is multi-task learning
- **Motivation**
  - Why we want to group tasks
- **Algorithms**
  - How to discover the grouping
- **Empirical results**
  - Validate our approach
- **Conclusion**
  - Summary
  - Future work

# Motivation

Existing work on multi-task feature learning

- ***single regularization term***
- ***All*** tasks are related.

$$\min_{W} \quad \sum_{t=1}^{T} loss(w_t, X_t, Y_t) + \underline{\lambda \|W\|_{tr}^2}$$



matrix W

[ Argyriou, et al. 2008. ]

# Motivation

When models are in ***mixture*** of subspaces

$$\min_{W} \quad \sum_{t=1}^{T} loss(w_t, X_t, Y_t) + \lambda \|W\|_{tr}^2$$

- Suboptimal to force with one regularizer
- Ex: synthetic data (later in the talk)



matrix W

# Motivation

When groups are given

$$\min_{W_1, W_2} \sum_{t=1}^{T} loss(w_t, X_t, Y_t) + \lambda \|W_1\|_{tr}^2 + \lambda \|W_2\|_{tr}^2$$

**Desiderata**

Regularize each group *separately*.

Automatically learn with whom to share



matrix $W_1$

matrix W

matrix $W_2$

# Outline

# Step1: use indicator matrix

Reformulate with task group assignment matrix **Q**



$$\|W_1\|_{tr}^2 = \|Q_1 W\|_{tr}^2$$

# Integer programming for Inferring with whom to share

*Re-formulate* with matrix Q
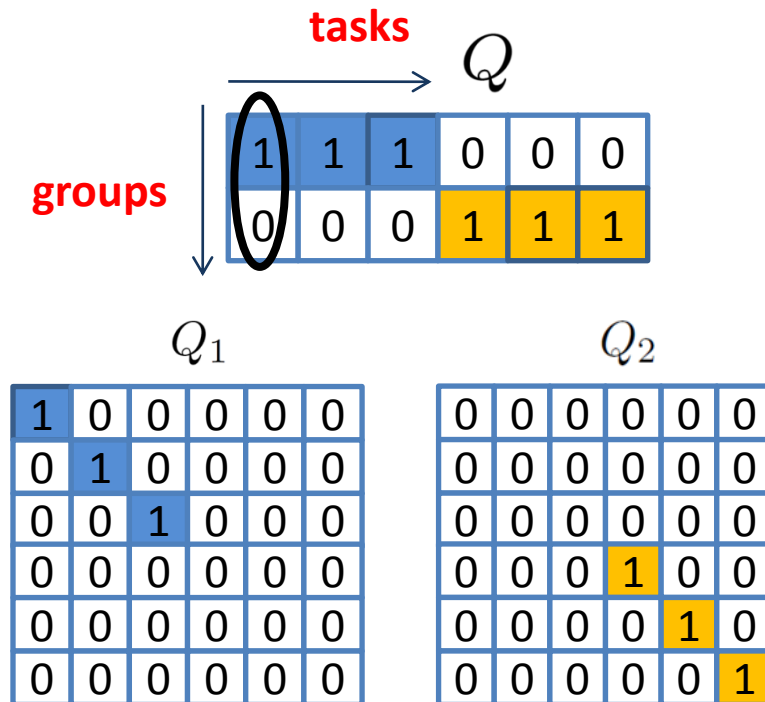
- *Integer* constraint
- *Hard* group assignment

$$\min_{W,Q} \sum_{t=1}^{T} loss(w_t, X_t, Y_t)$$

$$+ \lambda \|Q_1 W\|_{tr}^2 + \lambda \|Q_2 W\|_{tr}^2$$



$$\text{s.t} \quad q_{gt} \in \{0, 1\}$$

$$Q_1 + Q_2 = I$$

# Step 2: relax the constraint

- Approach 1:

  convex relaxation

  - *Continuous* constraint
  - Convex but *fractional* solutions

- Approach 2:

  non-convex relaxation

  - Use *square root of Q:* non-convex but *integer* solutions

$$\min_{W,Q} \sum_{t=1}^{T} loss(w_t, X_t, Y_t)$$

$$+ \lambda \|\sqrt{Q_1}W\|_{tr}^2 + \lambda \|\sqrt{Q_2}W\|_{tr}^2$$

$$\text{s.t} \quad 0 \le q_{gt} \le 1$$

$$Q_1 + Q_2 = I$$

# Integer solutions guaranteed

**Theorem 1.** *Let $\{Q_g^*\}$ be either the solution or a local optimum to the following optimization,*

$$\min \quad T(Q) = \sum_g \|W\sqrt{Q_g}\|_*^2$$

$$\text{s.t} \quad \sum_g Q_g = I \text{ with } 0 \leq q_{gt} \leq 1 \tag{9}$$

*then either one of the following is true: i) $\{Q_g^*\}$ is binary; ii) there exists another binary $\{Q_g'\}$ such that $T(Q^*) = T(Q')$.*

# Proofs.

$$T(\boldsymbol{Q}) = \sum_g \min_{\boldsymbol{\Omega}_g} \mathsf{Trace} \left[ \boldsymbol{\Omega}_g^{-1} \boldsymbol{W} \sqrt{\boldsymbol{Q}_g} \sqrt{\boldsymbol{Q}_g}^{\mathrm{T}} \boldsymbol{W}^{\mathrm{T}} \right] \quad (1)$$

where $\boldsymbol{\Omega}_g$ is constrained to be positive definitive. Furthermore, $\mathsf{Trace}[\boldsymbol{\Omega}_g] = 1$. Let $\boldsymbol{\Psi}_g = \boldsymbol{W}^{\mathrm{T}} \boldsymbol{\Omega}_g^{-1} \boldsymbol{W}$, we have

$$T(\boldsymbol{Q}) = \min \sum_g \mathsf{Trace} \left[ \boldsymbol{\Psi}_g \boldsymbol{Q}_g \right] \quad (2)$$

Since $\boldsymbol{Q}_g$ is a diagonal matrix, we have immediately

$$T(\boldsymbol{Q}) = \min \sum_g \sum_t \psi_{tt}^g q_{gt} \quad (3)$$

# Numerical Optimization

Optimize W and Q *iteratively*

- Fix Q, update W
  - *For each group, we solve*

$$\min \sum_{t:q_{gt}=1} \ell(\mathcal{D}_t; \boldsymbol{w}_t) + \gamma \|\boldsymbol{W}_g\|_*^2$$

  - *Use existing algorithm*

    *cf: Argyriou, et al. **Convex multi-task feature learning**. MLJ 2008.*

# **Numerical Optimization**

Optimize W and Q *iteratively*

- Fix W, update Q

  - *Use gradient descent*

$$
\min_Q \quad \sum_g \| \sqrt{Q_g} W \|_{tr}^2
$$

$$
\text{s.t} \quad \sum_g Q_g = I \text{ with } 0 \le q_{gt} \le 1
$$

  - *Remove constraints*

    - *by re-parameterization: α is unconstrained*

$$
q_{gt} = \frac{e^{\alpha_{gt}}}{\sum_{g=1}^{G} e^{\alpha_{gt}}} \quad \textbf{( soft assigning )}
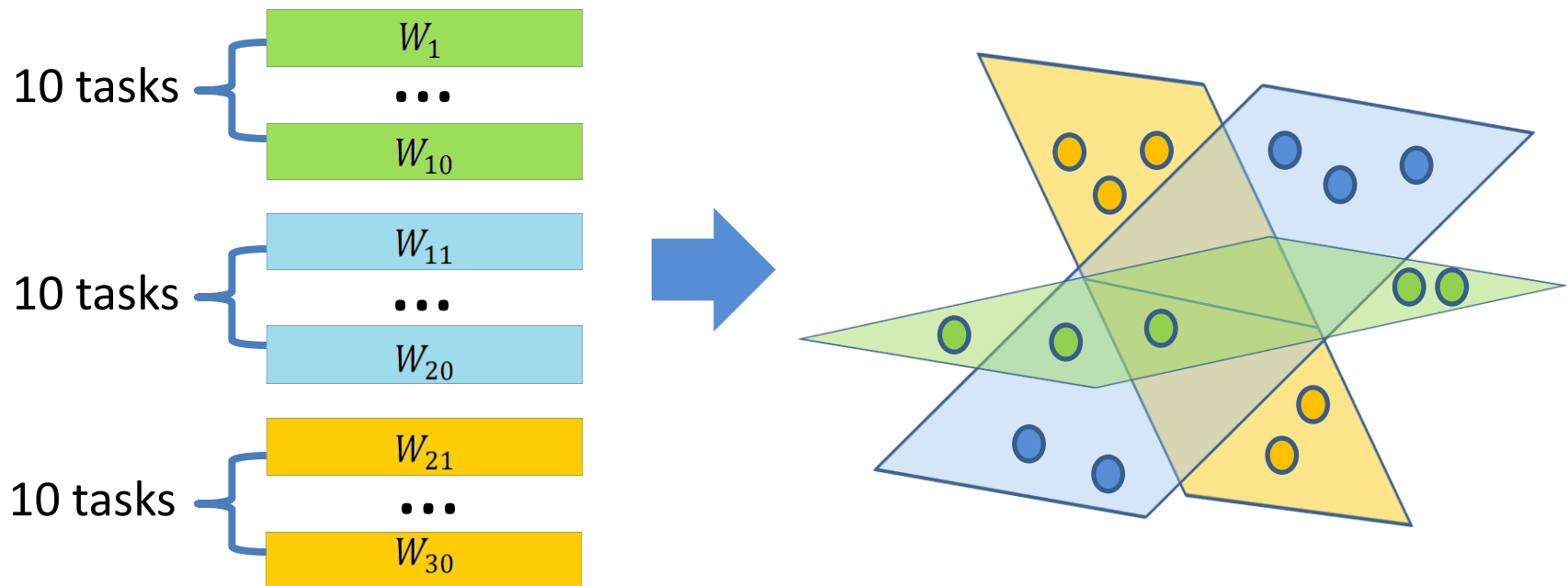$$

# Outline

- **Background**
  - What is multi-task learning
- **Motivation**
  - Why we want to group tasks
- **Algorithms**
  - How to discover the grouping
- **Empirical results**
  - Validate our approach
- **Conclusion**
  - Summary
  - Future work

# Results: synthetic data

## Setup

– We have 30 tasks with 3 groups (10 tasks per group).

– Each task is a regression problem.
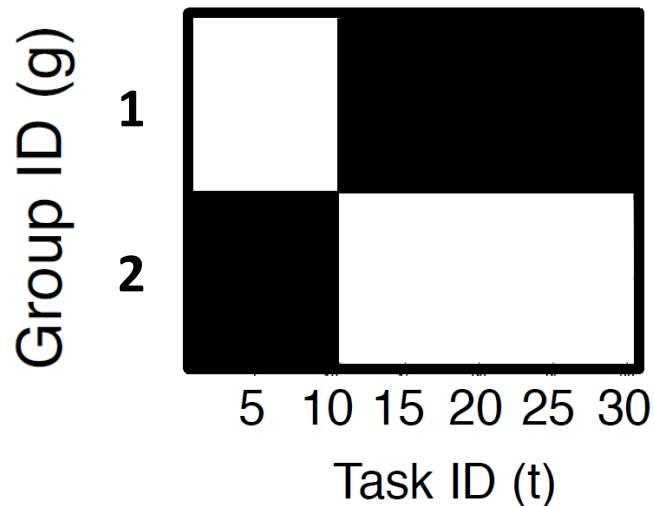
– Tasks in the same group *use the same feature.*

# Grouping results of the tasks

- Specify the correct number of groups
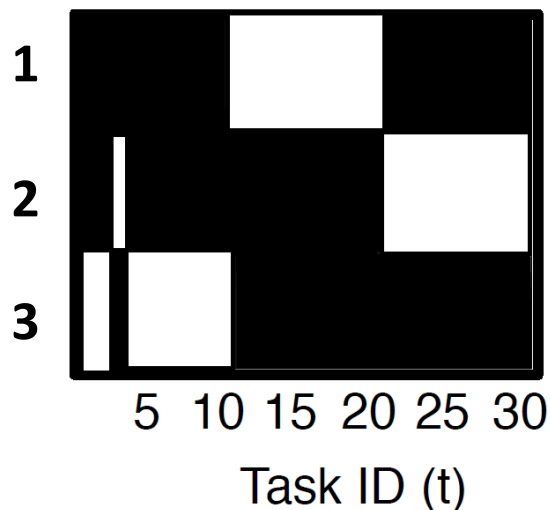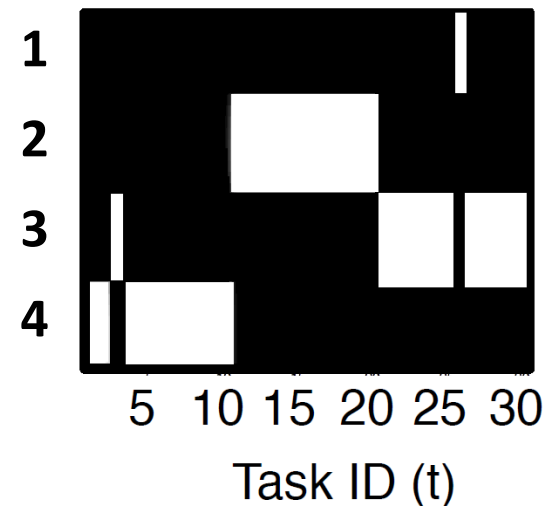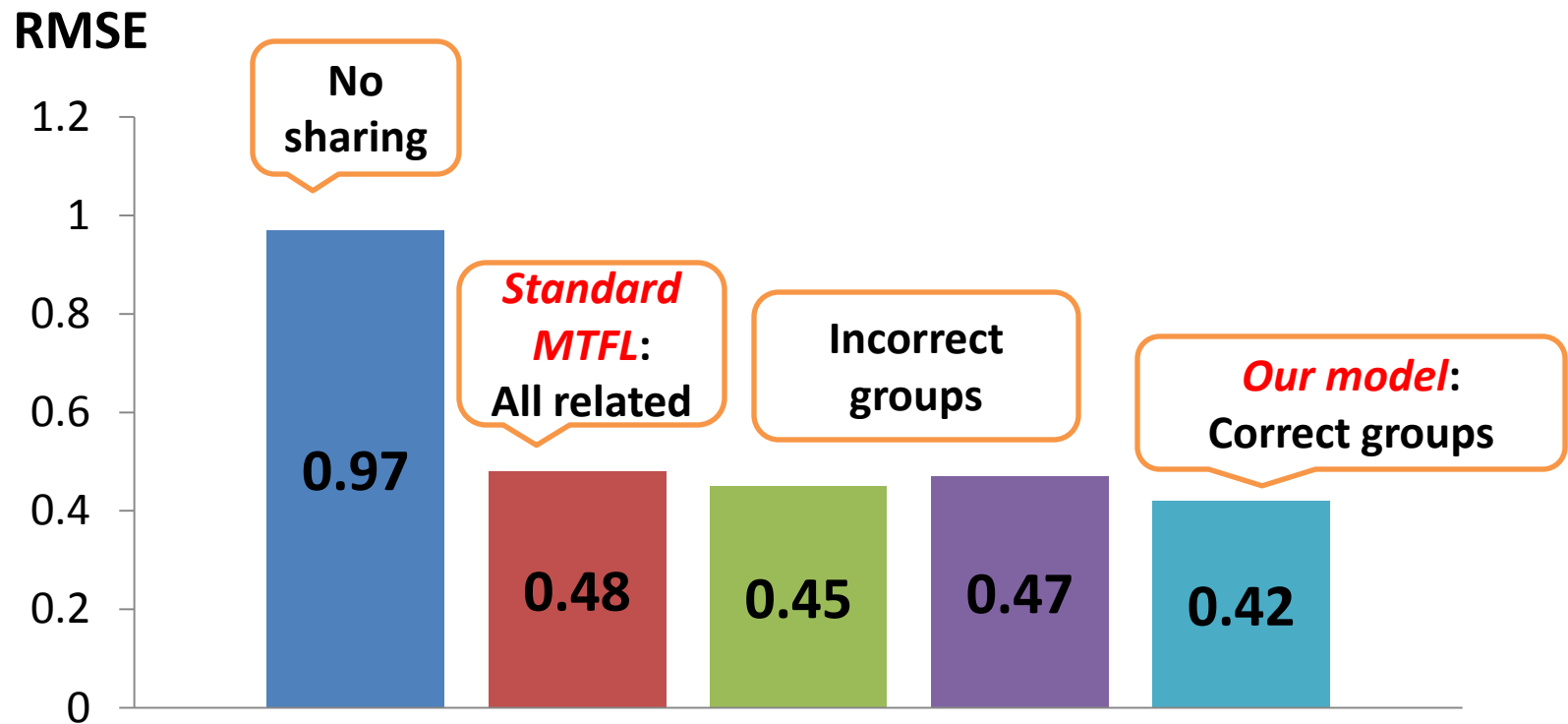- Identify the **correct grouping**

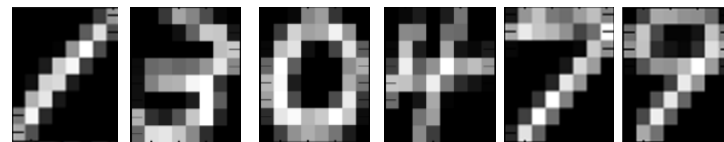# Also improve generalization

- Measure **average *root-mean-square error***.

- Obtain best performance with ***correct grouping***.

RMSE

**No sharing**

***Standard MTFL*:** **All related**

**Incorrect groups**

***Our model*:** **Correct groups**

1.2

1
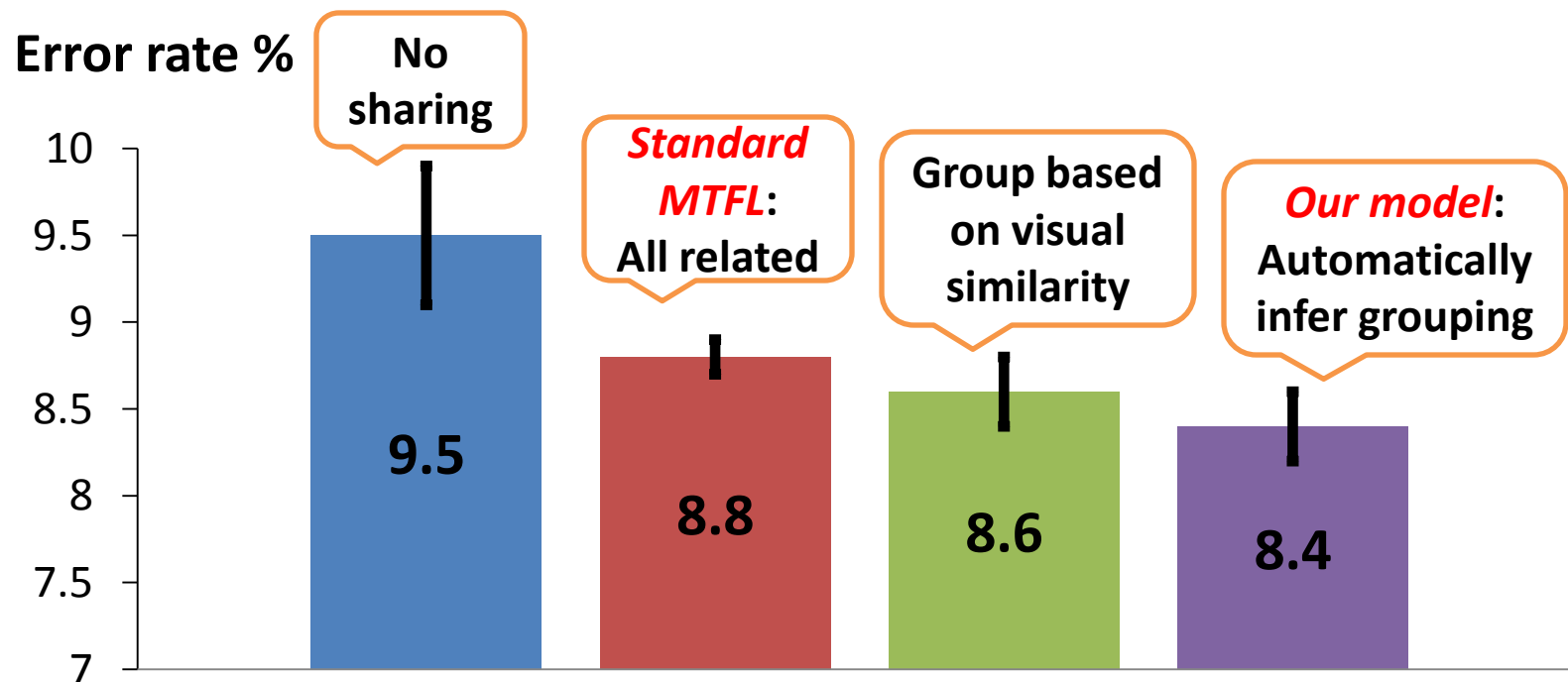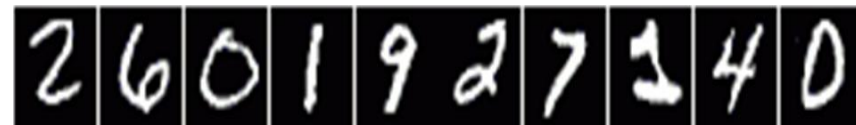
0.8

0.6

0.4

0.2

0

0.97

0.48

0.45

0.47

0.42

# Results: USPS

## Setup

– 10-way classification on images of 10 handwritten digits

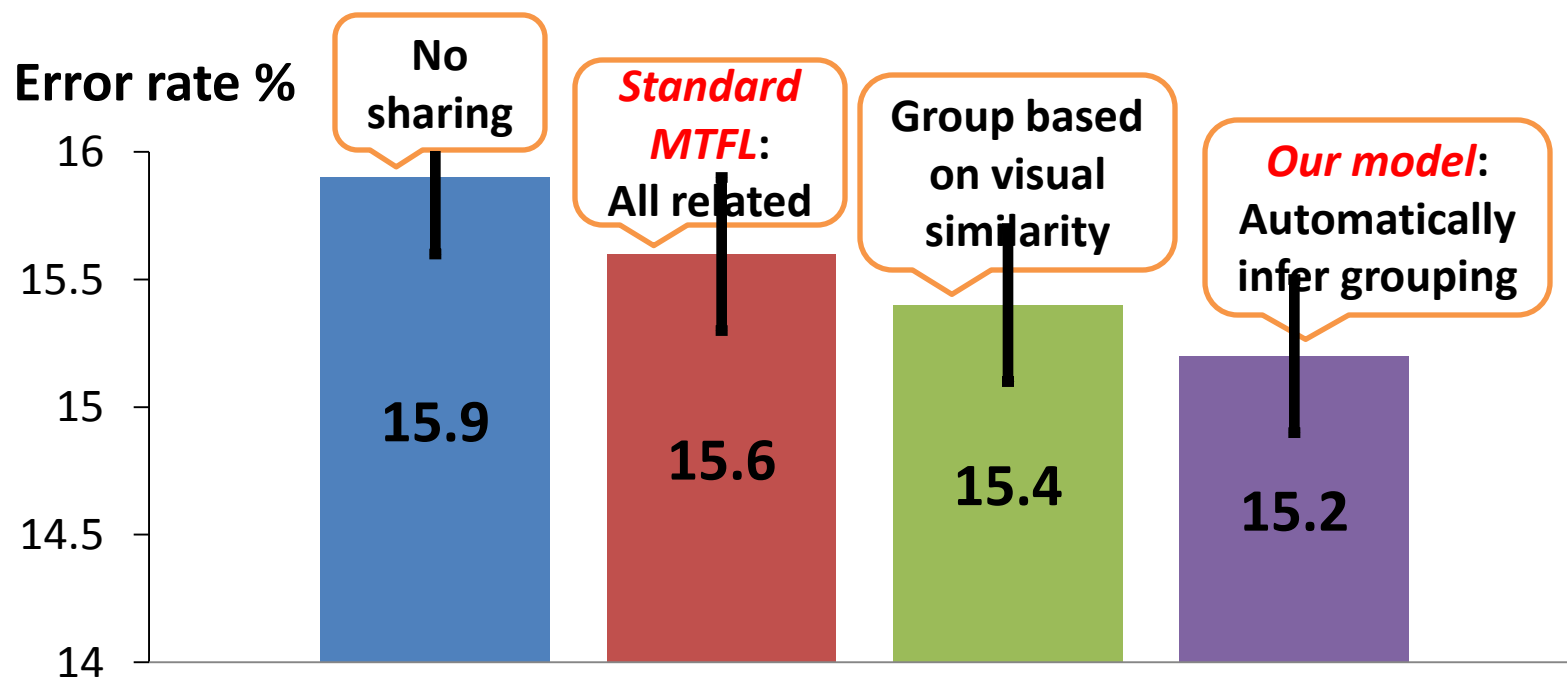– 1000 training data

– Classifier: binary logistic regression

**Error rate %**

| | No sharing | Standard MTFL: All related | Group based on visual similarity | Our model: Automatically infer grouping |
|---|---|---|---|---|
| Error rate | 9.5 | 8.8 | 8.6 | 8.4 |

# Results: MNIST

## Setup

– 10-way classification on images of 10 handwritten digits

– 1000 training data

– Classifier: binary logistic regression
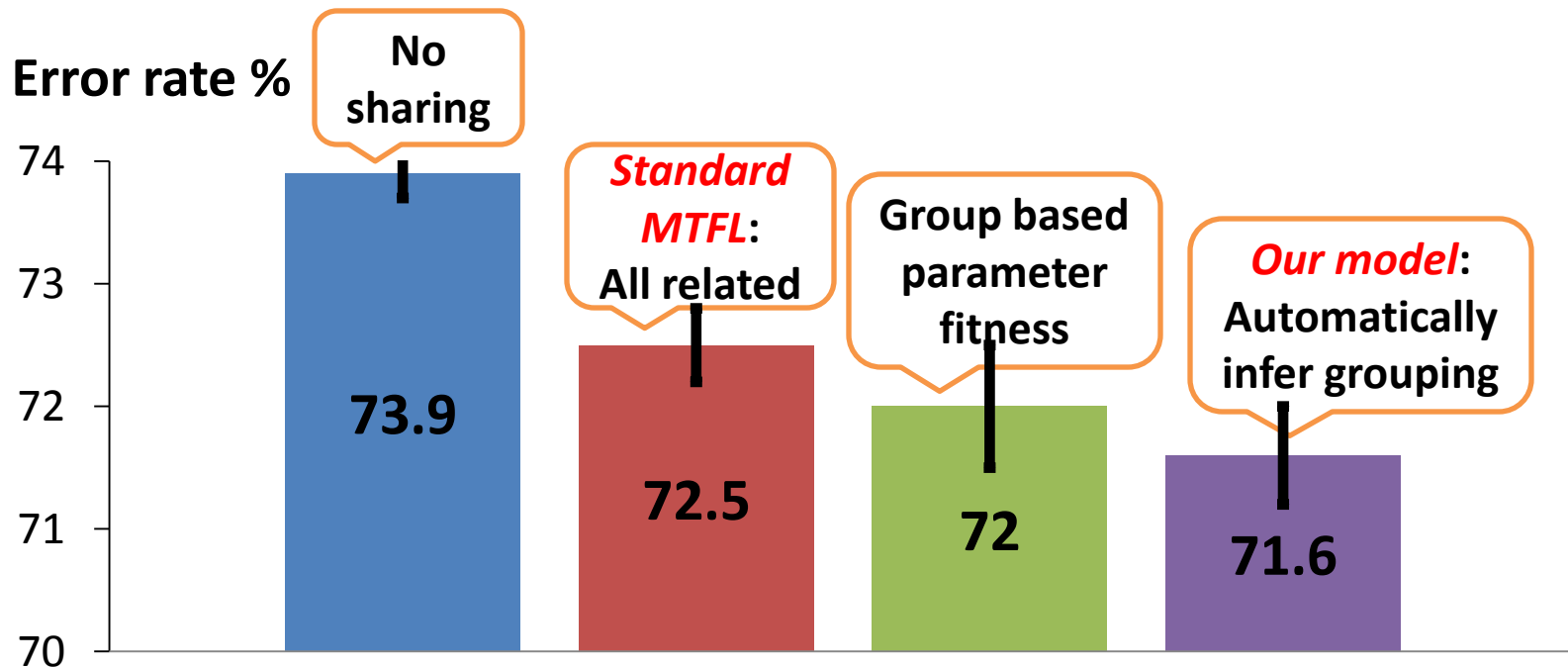
**Error rate %**

Bar chart showing error rates:
- **No sharing**: 15.9
- **Standard MTFL: All related**: 15.6
- **Group based on visual similarity**: 15.4
- **Our model: Automatically infer grouping**: 15.2

# Results: recognize animals

## Setup

– Data set: Animal with Attributes (images of 20 classes)

– 1000 training data; Features: SIFT

– Classifier: binary logistic regression

# Grouping results on digits data

# Grouping results on animal data
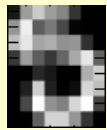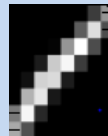
**Animal with Attributes data set**: 20 classes are used

# Comparison with other methods



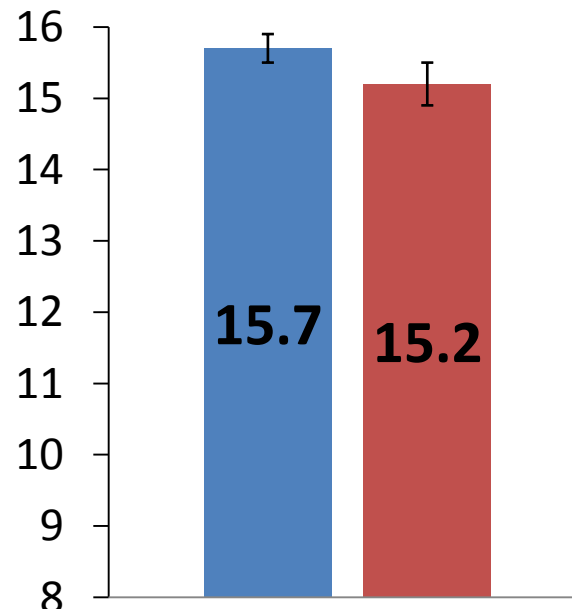■ **Online learning of sets of kernels** [Argyiou, et al. ECML 2008]
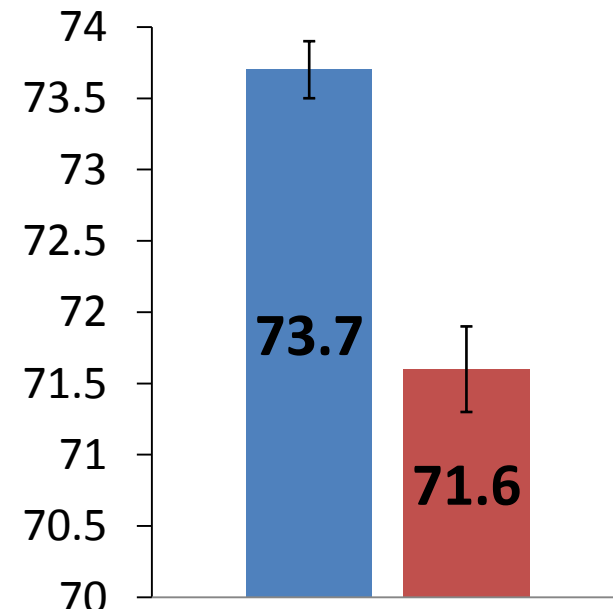■ **Our method**

# Conclusions

Multi-task feature learning

- Beneficial to **identify related tasks.**

- **Sharing with related tasks** instead all of them helps.

- Effective joint inference of **shared structure** and **model parameters**.

Future work

- More complex structures.

- Investigation of the grouping robustness.

- Transfer for new tasks.

# Transfer for new tasks

- Q and W of old tasks is enough
- keep Q and W of old tasks fixed
- update Q and W of new tasks

new Q          old Q (fixed)

| ? | ? | 1 | 1 | 1 | 0 | 0 | 0 |
| ? | ? | 0 | 0 | 0 | 1 | 1 | 1 |
| ? | ? | 0 | 0 | 0 | 0 | 0 | 0 |
| ? | ? | 0 | 0 | 0 | 0 | 0 | 0 |

$w_{new}$
$w_{new}$

new W

$w_1$
$w_2$
$w_3$
$w_4$
$w_5$
$w_6$

old W (fixed)