# Ten years of pedestrian detection, what have we learned ?

Rodrigo Benenson

Mohamed Omran

Jan Hosang

Markus Mathias

Shanshan Zhang

# This presentation:
# what works and does not work for pedestrian detection ?

[Benenson et al. ECCVw 2014]

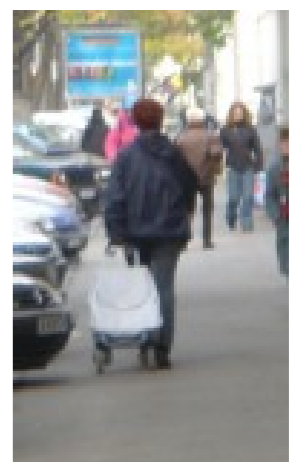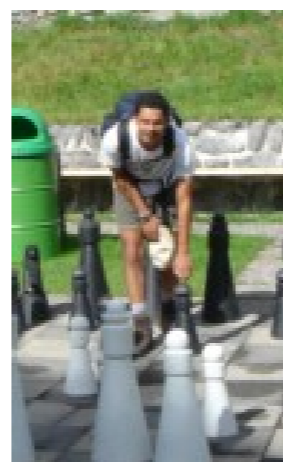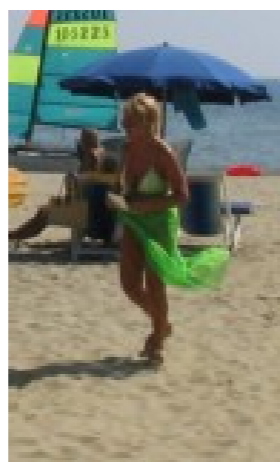"Science is the belief in the ignorance of experts"

Richard Feynman

# Why pedestrian detection ?

Rodrigo Benenson |  CMP Colloquium 2014

# Pedestrian detection is an interesting problem

- Large variance for intra-class appearance

- Strong illumination changes

- Deformations

- Occlusions

- (Interest on small instances)

- No structural variations
(number of wings in an airplane)

# Pedestrian detection is harder than you might think



INRIA training examples

# Pedestrian detection is harder than you might think



To be or not to be pedestrian ? (Caltech test set)

# Pedestrian detection is harder than you might think



To be or not to be pedestrian ?  (Caltech test set)

# Pedestrian detection is mature

1) Many ideas have been proposed
   ⇨1000+ papers with "pedestrian detection" title

2)

3)

# Pedestrian detection is mature

1) Many ideas have been proposed

2) Good enough benchmarks are available

3)



Caltech-USA dataset

KITTI dataset

[Dollar, Wojek, Schiele, Perona 2009]
[Geiger et al. 2013]

# Pedestrian detection is mature

1) Many ideas have been proposed

2) Good enough benchmarks are available

3) Well defined metric

⇒ Average miss-rate
(lower is better)

INRIA dataset



Better

miss rate

Shapelet-orig (90.5%)
PoseInvSvm (68.6%)
VJ-OpenCv (53.0%)
PoseInv (51.4%)
Shapelet (50.4%)
VJ (47.5%)
FtrMine (34.0%)
Pls (23.4%)
HOG (23.1%)
HikSvm (21.9%)
LatSvm-V1 (17.5%)
MultiFtr (15.6%)
MultiFtr+CSS (10.9%)
LatSvm-V2 (9.3%)
FPDW (9.3%)
ChnFtrs (8.7%)

false positives per image

Better

# Pedestrian detection is mature, but not stagnant



Papers with "Pedestrian detection" in the title

# Great progress in pedestrian detection during last decade



## Caltech-USA is currently the most active dataset.

# Many different ideas have been explored

- Sophisticated features

- Deformable parts

- Deeper architectures

- Non-linear classifiers

- Richer training data

- Geometric priors

- Motion information



INRIA dataset

Better

miss rate

Shapelet–orig (90.5%)
PoseInvSvm (68.6%)
VJ–OpenCv (53.0%)
PoseInv (51.4%)
Shapelet (50.4%)
VJ (47.5%)
FtrMine (34.0%)
Pls (23.4%)
HOG (23.1%)
HikSvm (21.9%)
LatSvm–V1 (17.5%)
MultiFtr (15.6%)
MultiFtr+CSS (10.9%)
LatSvm–V2 (9.3%)
FPDW (9.3%)
ChnFtrs (8.7%)

Better

false positives per image

# More is more

# More is more

## Less is more

# Less is more

- ~~Sophisticated features~~
- ~~Deformable parts~~
- ~~Deeper architectures~~
- ~~Non-linear classifiers~~
- ~~Richer training data~~
- ~~Geometric priors~~
- ~~Motion information~~



INRIA dataset

Better

Better

miss rate

false positives per image

81.08% Shapelet
79.81% PoseInv
72.05% VJ
57.70% FtrMine
45.18% HOG
43.53% LatSvm-V1
41.16% HikSvm
38.69% Pls
37.25% HogLbp
35.40% MultiFtr
30.91% FeatSynth
23.93% MultiFtr+CSS
23.49% MLS
20.53% FPDW
20.44% ChnFtrs
19.55% LatSvm-V2
18.92% EBLearn
18.26% CrossTalk
18.21% Ours-SquaresChnFtrs
15.40% VeryFast
13.06% Ours-Roerei

[Benenson et al. CVPR 2013]

# Revisiting the basics:
# what makes pedestrian detection <u>really</u> work ?

[Benenson et al. CVPR 2013]

Message of the day:

One (simple and effective) core

+

3 add-ons

# Quick chronology: 5 landmarks



**40+ methods currently on Caltech-USA**

# Quick chronology: 5 landmarks



[Viola & Jones 2004]

94.73% VJ

miss rate

false positives per image

# Quick chronology: 5 landmarks



[Viola & Jones 2004]

[Dalal & Triggs 2005]

94.73% VJ

68.46% HOG

# Quick chronology: 5 landmarks



[Viola & Jones 2004]
[Dalal & Triggs 2005]

94.73% VJ

68.46% HOG

63.26% LatSvm−V2

[Felzenszwalb et al. 2008]
**DPM family**

# Quick chronology: 5 landmarks



[Viola & Jones 2004]
[Dalal & Triggs 2005]
[Felzenszwalb et al. 2008]

[Dollar et al. 2009]
**DF family**

94.73% VJ

68.46% HOG

63.26% LatSvm−V2

56.34% ChnFtrs

# Quick chronology: 5 landmarks



[Viola & Jones 2004]
[Dalal & Triggs 2005]
[Felzenszwalb et al. 2008]
[Dollar et al. 2009]

94.73% VJ
68.46% HOG
63.26% LatSvm−V2
56.34% ChnFtrs
53.14% DBN−Isol

[Ouyang & Wang 2012]
**DN family**

# Quick chronology: 5 landmarks



[Viola & Jones 2004]
[Dalal & Triggs 2005]
[Felzenszwalb et al. 2008]
[Dollar et al. 2009]
[Ouyang & Wang 2012]

94.73% VJ
68.46% HOG
63.26% LatSvm−V2
56.34% ChnFtrs
53.14% DBN−Isol
22.49% Ours−Katamari

miss rate

false positives per image

[Benenson et al. 2014]

# Quick chronology: 5 landmarks



[Viola & Jones 2004]
[Dalal & Triggs 2005]
[Felzenszwalb et al. 2008]
[Dollar et al. 2009]
[Ouyang & Wang 2012]
[Benenson 2014]

94.73% VJ
68.46% HOG
63.26% LatSvm−V2
56.34% ChnFtrs
53.14% DBN−Isol
22.49% Ours−Katamari

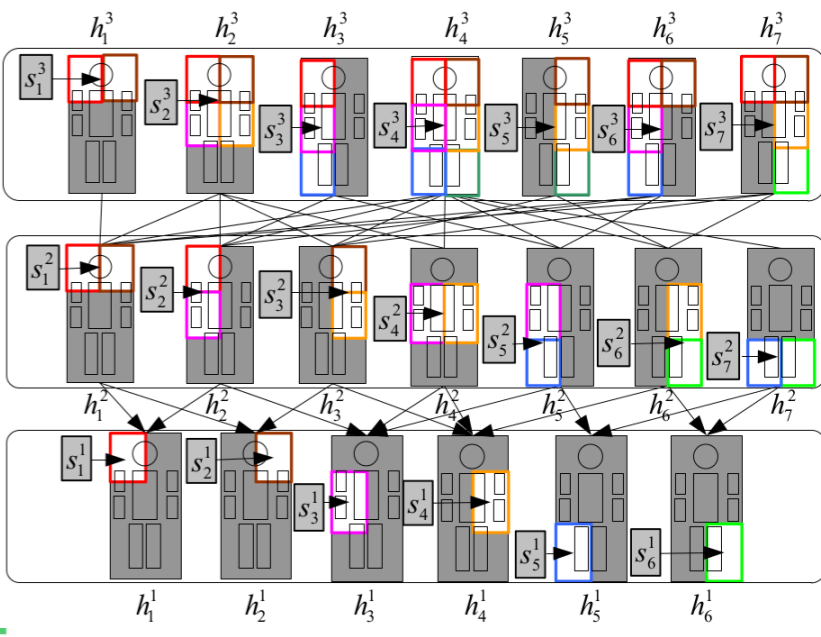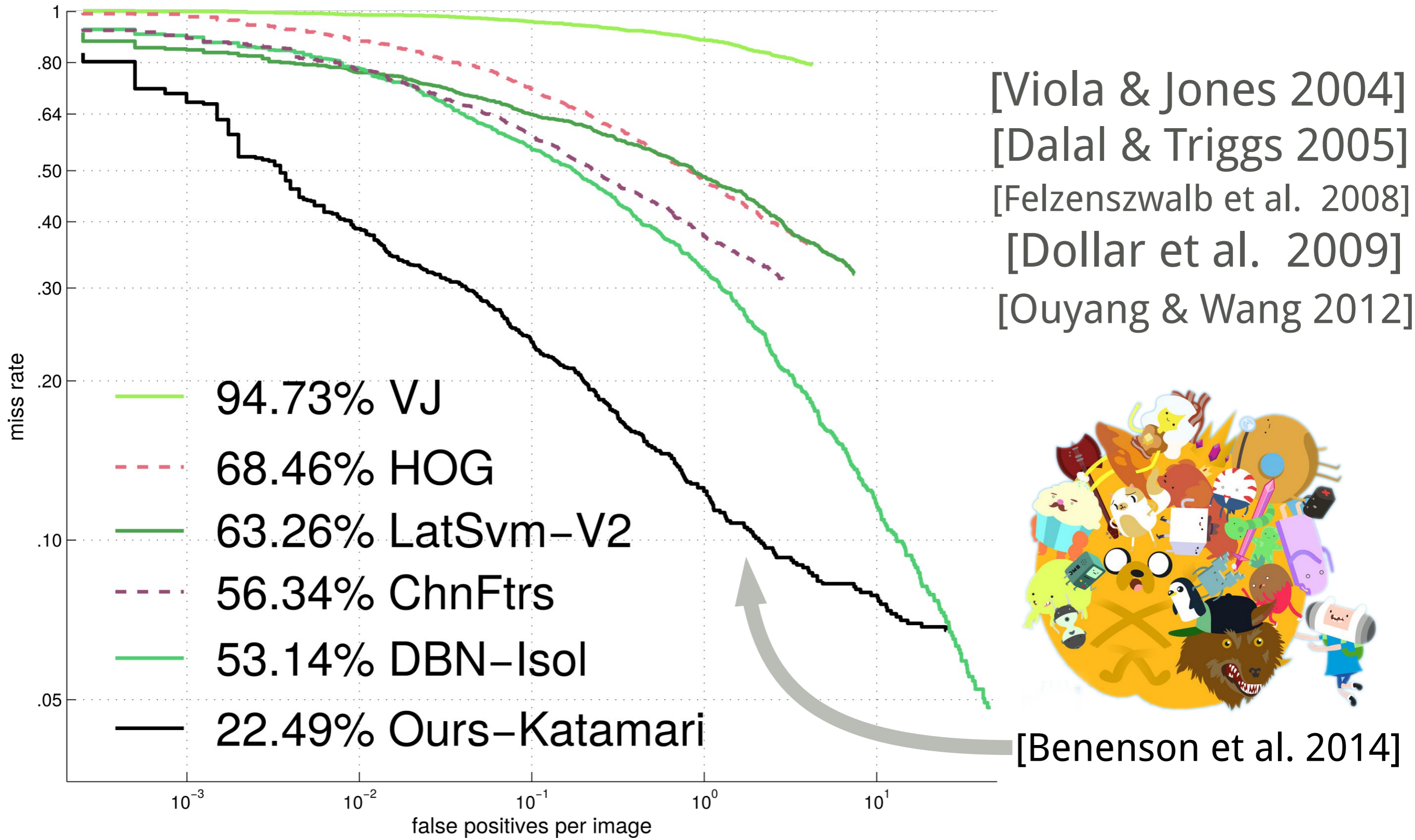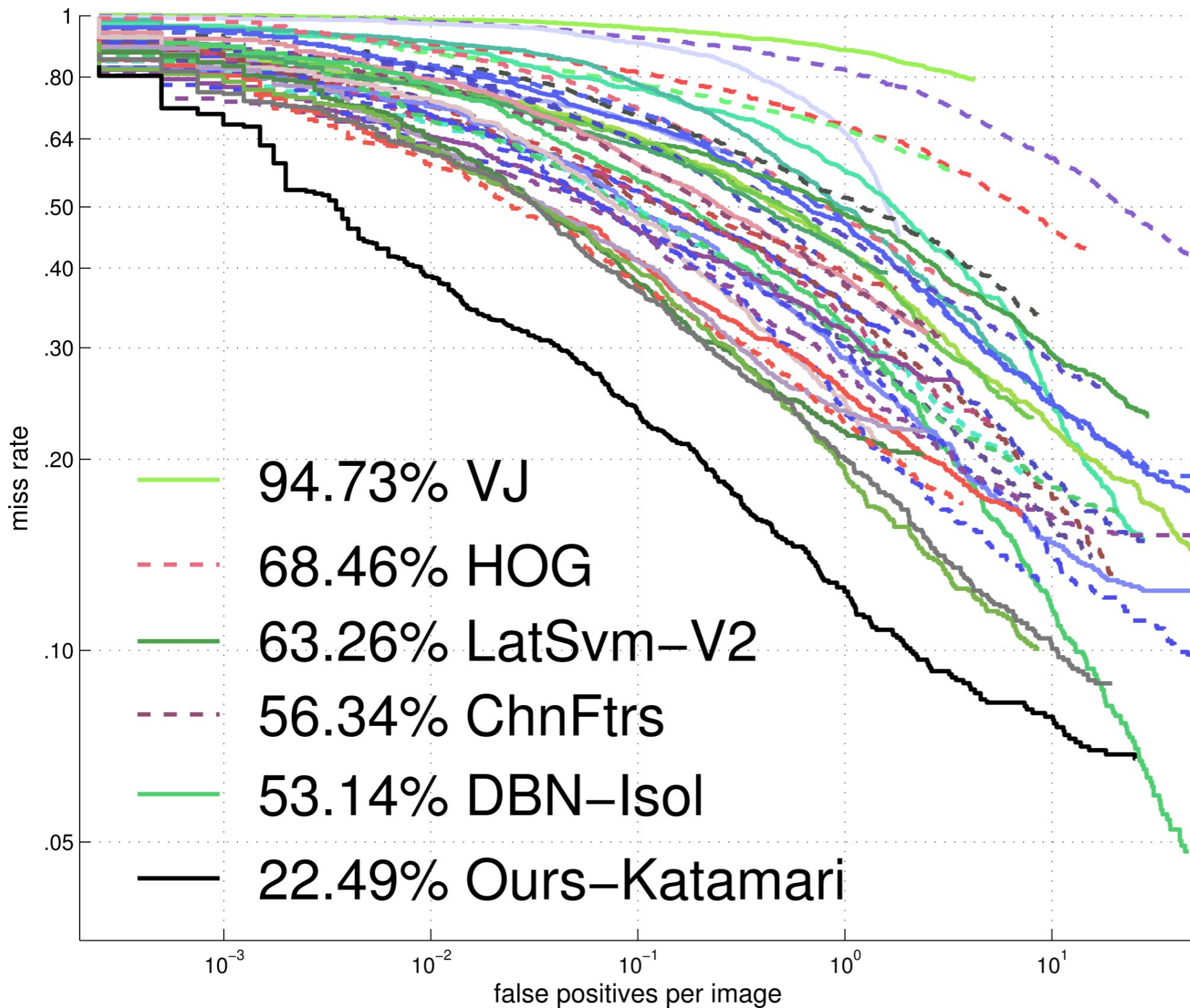| Method MR | Family | Features | Classifier | Context | Deep | Parts | M-Scales | More data | Feat. type | Training |
|---|---|---|---|---|---|---|---|---|---|---|
| VJ [9] 94.73% | DF | ✓ | ✓ | | | | | | Haar | I |
| Shapelet [10] 91.37% | - | ✓ | | | | | | | Gradients | I |
| PoseInv [11] 86.32% | - | | | | | ✓ | | | HOG | I+ |
| LatSvm-V1 [12] 79.78% | DPM | | | | | ✓ | | | HOG | P |
| ConvNet [13] 77.20% | DN | | | | ✓ | | | | Pixels | I |
| FtrMine [14] 74.42% | DF | ✓ | | | | | | | HOG+Color | I |
| HikSvm [15] 73.39% | - | | ✓ | | | | | | HOG | I |
| HOG [1] 68.46% | - | ✓ | ✓ | | | | | | HOG | I |
| MultiFtr [16] 68.26% | DF | ✓ | ✓ | | | | | | HOG+Haar | I |
| HogLbp [17] 67.77% | - | ✓ | | | | | | | HOG+LBP | I |
| AFS+Geo [18] 66.76% | - | | | ✓ | | | | | Custom | I |
| AFS [18] 65.38% | - | | | | | | | | Custom | I |
| LatSvm-V2 [19] 63.26% | DPM | | ✓ | | | ✓ | | | HOG | I |
| Pls [20] 62.10% | - | ✓ | ✓ | | | | | | Custom | I |
| MLS [21] 61.03% | DF | ✓ | | | | | | | HOG | I |
| MultiFtr+CSS [22] 60.89% | DF | ✓ | | | | | | | Many | T |
| FeatSynth [23] 60.16% | - | ✓ | ✓ | | | | | | Custom | I |
| pAUCBoost [24] 59.66% | DF | ✓ | ✓ | | | | | | HOG+COV | I |
| FPDW [25] 57.40% | DF | | | | | | | | HOG+LUV | I |
| ChnFtrs [26] 56.34% | DF | ✓ | ✓ | | | | | | HOG+LUV | I |
| CrossTalk [27] 53.88% | DF | | | ✓ | | | | | HOG+LUV | I |
| DBN-Isol [28] 53.14% | DN | | | | ✓ | | | | HOG | I |
| ACF [29] 51.36% | DF | ✓ | | | | | | | HOG+LUV | I |
| RandForest [30] 51.17% | DF | | ✓ | | | | | | HOG+LBP | I&C |
| MultiFtr+Motion [22] 50.88% | DF | ✓ | | | | | | ✓ | Many+Flow | T |
| *SquaresChnFtrs* [31] 50.17% | DF | ✓ | | | | | | | HOG+LUV | I |
| Franken [32] 48.68% | DF | | ✓ | | | | | | HOG+LUV | I |
| MultiResC [33] 48.45% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C |
| Roerei [31] 48.35% | DF | ✓ | | | | | ✓ | | HOG+LUV | I |
| DBN-Mut [34] 48.22% | DN | | | ✓ | ✓ | | | | HOG | C |
| MF+Motion+2Ped [35] 46.44% | DF | | | ✓ | | | | ✓ | Many+Flow | I+ |
| MOCO [36] 45.53% | - | ✓ | | ✓ | | | | | HOG+LBP | C |
| MultiSDP [37] 45.39% | DN | ✓ | | ✓ | ✓ | | | | HOG+CSS | C |
| ACF-Caltech [29] 44.22% | DF | ✓ | | | | | | | HOG+LUV | C |
| MultiResC+2Ped [35] 43.42% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C+ |
| WordChannels [38] 42.30% | DF | ✓ | | | | | | | Many | C |
| MT-DPM [39] 40.54% | DPM | | | | | ✓ | ✓ | | HOG | C |
| JointDeep [40] 39.32% | DN | | | ✓ | | | | | Color+Gradient | C |
| SDN [41] 37.87% | DN | | | | ✓ | ✓ | | | Pixels | C |
| MT-DPM+Context [39] 37.64% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C+ |
| ACF+SDt [42] 37.34% | DF | ✓ | | | | | | ✓ | ACF+Flow | C+ |
| *SquaresChnFtrs* [31] 34.81% | DF | ✓ | | | | | | | HOG+LUV | C |
| InformedHaar [43] 34.60% | DF | ✓ | | | | | | | HOG+LUV | C |
| *Katamari-v1* 22.49% | DF | ✓ | | ✓ | | | | ✓ | HOG+Flow | C+ |

# What is driving the quality progress ?

- solution family (DPM, deep networks, decision forests)

- better classifiers

- deformable parts

- multi-scale models

- deep architectures

- training data

- additional (test time) data

- exploiting context

- better features

# What is driving the quality progress ?

- solution family (DPM, deep networks, decision forests)

- better classifiers

- deformable parts

- multi-scale models

- deep architectures

- training data

- additional (test time) data

- exploiting context

- better features

# What is driving the quality progress ?

- **solution family
  (DPM, deep networks, decision forests)**

- **better classifiers**

- deformable parts

- multi-scale models

- deep architectures

- training data

- additional (test time) data

- exploiting context

- better features

# Surprise 1:

There is no clear winner regarding solution family (DPM, DN, or DF) or classifier type.


SURPRISE!

# What is driving the quality progress ?

- ~~solution family (DPM, deep networks, decision forests)~~

- ~~better classifiers~~

- deformable parts

- multi-scale models

- deep architectures

- training data

- additional (test time) data

- exploiting context

- better features

# Data is inconclusive: the DPM case

Latent-SVM v2 ⇒ 63%
[Felzenszwalb et al.  2010]

MultiResC ⇒ 49%
[Park et al.  2010]

MT-DPM ⇒ 41%
[Yan et al.  2013]

?

Vanilla DPM v4 ⇒ **42%**
[Yan et al.  2014]

!

Our rigid template ⇒ 34%
[Benenson  2014]

[Hariharan et al. CVPR 2014]
[Girshick et al. arXiv 2014]

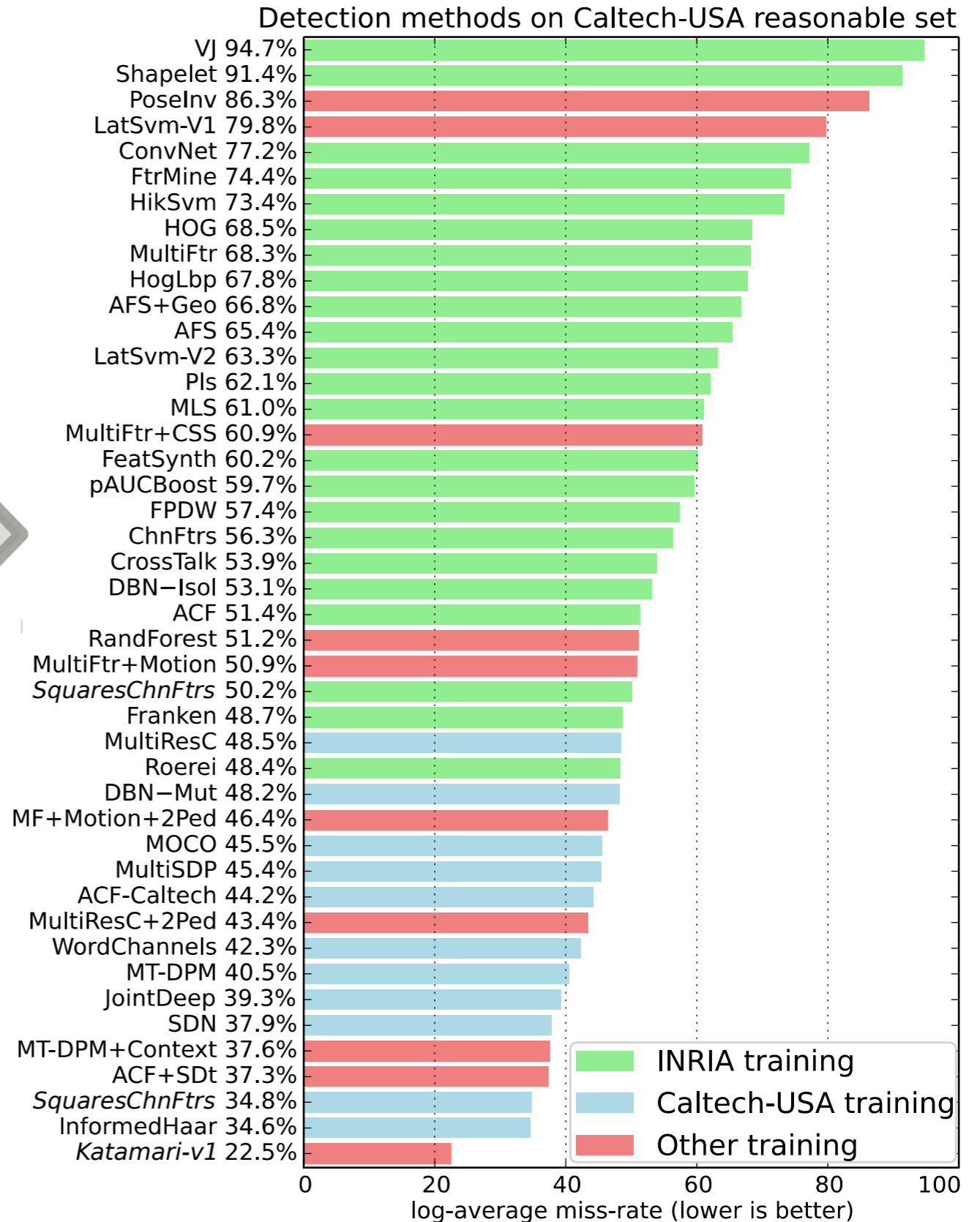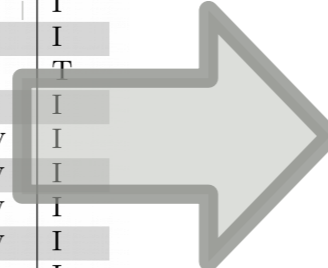# What is driving the quality progress ?

- ~~solution family (DPM, deep networks, decision forests)~~

- ~~better classifiers~~

- ~~deformable parts~~

- ~~multi-scale models~~

- ~~deep architectures~~

- training data

- additional (test time) data

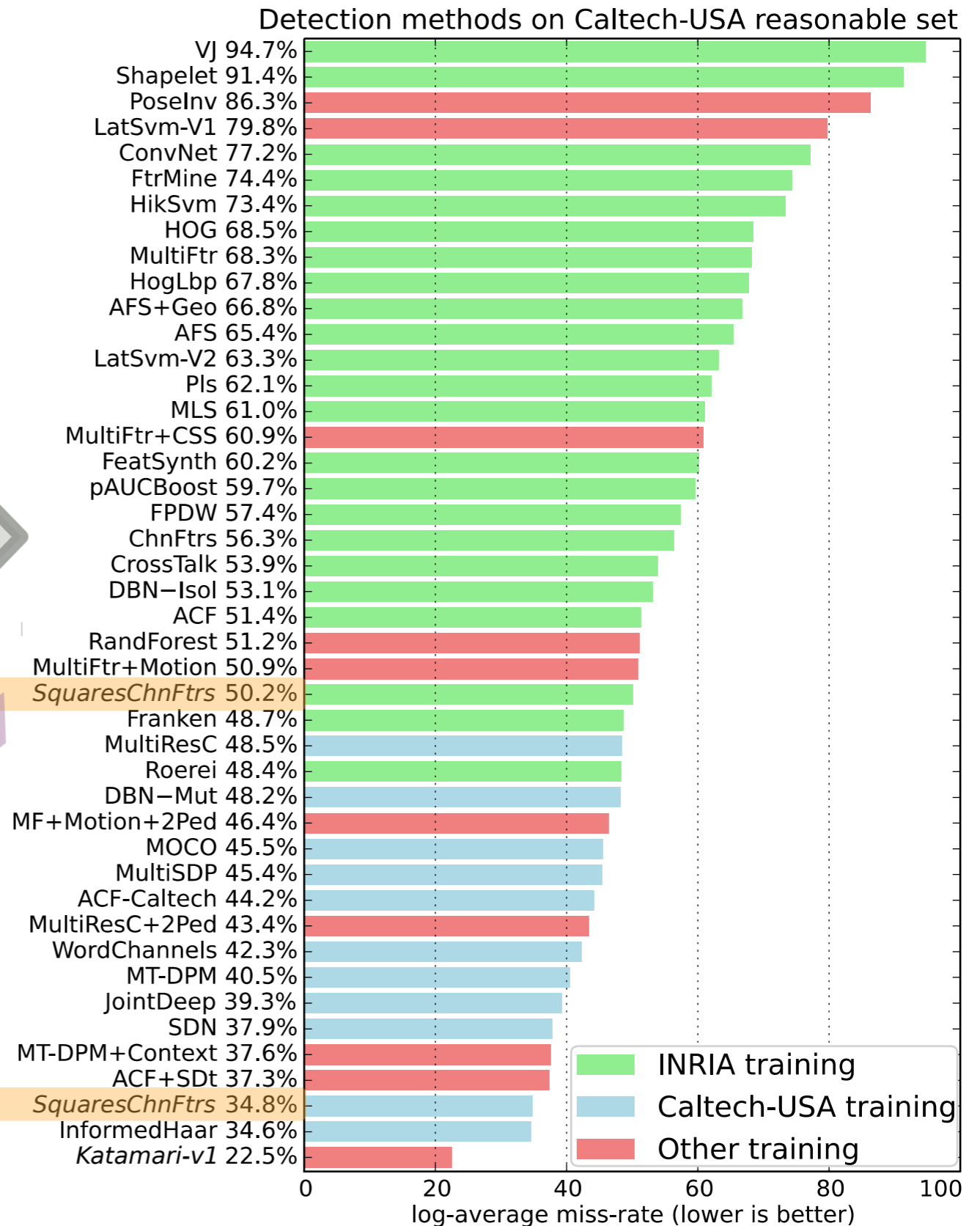- exploiting context

- better features

# Training data matters (you knew this already)

| Method | MR | Family | Features | Classifier | Context | Deep | Parts | M-Scales | More data | Feat. type | Training |
|---|---|---|---|---|---|---|---|---|---|---|---|
| VJ [9] | 94.73% | DF | ✓ | ✓ | | | | | | Haar | I |
| Shapelet [10] | 91.37% | - | ✓ | | | | | | | Gradients | I |
| PoseInv [11] | 86.32% | - | | | ✓ | | | | | HOG | I+ |
| LatSvm-V1 [12] | 79.78% | DPM | | | ✓ | | | | | HOG | P |
| ConvNet [13] | 77.20% | DN | | | | ✓ | | | | Pixels | I |
| FtrMine [14] | 74.42% | DF | ✓ | | | | | | | HOG+Color | I |
| HikSvm [15] | 73.39% | - | | ✓ | | | | | | HOG | I |
| HOG [1] | 68.46% | - | ✓ | ✓ | | | | | | HOG | I |
| MultiFtr [16] | 68.26% | DF | ✓ | ✓ | | | | | | HOG+Haar | I |
| HogLbp [17] | 67.77% | - | ✓ | | | | | | | HOG+LBP | I |
| AFS+Geo [18] | 66.76% | - | | | ✓ | | | | | Custom | I |
| AFS [18] | 65.38% | - | | | | | | | | Custom | I |
| LatSvm-V2 [19] | 63.26% | DPM | ✓ | | ✓ | | | | | HOG | I |
| Pls [20] | 62.10% | - | ✓ | ✓ | | | | | | Custom | I |
| MLS [21] | 61.03% | DF | ✓ | | | | | | | HOG | I |
| MultiFtr+CSS [22] | 60.89% | DF | ✓ | | | | | | | Many | T |
| FeatSynth [23] | 60.16% | - | ✓ | ✓ | | | | | | Custom | I |
| pAUCBoost [24] | 59.66% | DF | ✓ | ✓ | | | | | | HOG+COV | I |
| FPDW [25] | 57.40% | DF | | | | | | | | HOG+LUV | I |
| ChnFtrs [26] | 56.34% | DF | ✓ | ✓ | | | | | | HOG+LUV | I |
| CrossTalk [27] | 53.88% | DF | | | ✓ | | | | | HOG+LUV | I |
| DBN-Isol [28] | 53.14% | DN | | | | ✓ | | | | HOG | I |
| ACF [29] | 51.36% | DF | ✓ | | | | | | | HOG+LUV | I |
| RandForest [30] | 51.17% | DF | | ✓ | | | | | | HOG+LBP | I&C |
| MultiFtr+Motion [22] | 50.88% | DF | ✓ | | | | | ✓ | | Many+Flow | T |
| SquaresChnFtrs [31] | 50.17% | DF | ✓ | | | | | | | HOG+LUV | I |
| Franken [32] | 48.68% | DF | | ✓ | | | | | | HOG+LUV | I |
| MultiResC [33] | 48.45% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C |
| Roerei [31] | 48.35% | DF | ✓ | | | | | ✓ | | HOG+LUV | I |
| DBN-Mut [34] | 48.22% | DN | | | ✓ | ✓ | | | | HOG | C |
| MF+Motion+2Ped [35] | 46.44% | DF | | | ✓ | | | ✓ | | Many+Flow | I+ |
| MOCO [36] | 45.53% | - | ✓ | | ✓ | | | | | HOG+LBP | C |
| MultiSDP [37] | 45.39% | DN | ✓ | | ✓ | ✓ | | | | HOG+CSS | C |
| ACF-Caltech [29] | 44.22% | DF | ✓ | | | | | | | HOG+LUV | C |
| MultiResC+2Ped [35] | 43.42% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C+ |
| WordChannels [38] | 42.30% | DF | ✓ | | | | | | | Many | C |
| MT-DPM [39] | 40.54% | DPM | | | | | ✓ | ✓ | | HOG | C |
| JointDeep [40] | 39.32% | DN | | | ✓ | | | | | Color+Gradient | C |
| SDN [41] | 37.87% | DN | | | | ✓ | ✓ | | | Pixels | C |
| MT-DPM+Context [39] | 37.64% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C+ |
| ACF+SDt [42] | 37.34% | DF | ✓ | | | | | | ✓ | ACF+Flow | C+ |
| SquaresChnFtrs [31] | 34.81% | DF | ✓ | | | | | | | HOG+LUV | C |
| InformedHaar [43] | 34.60% | DF | ✓ | | | | | | | HOG+LUV | C |
| Katamari-v1 | 22.49% | DF | ✓ | | ✓ | | | | ✓ | HOG+Flow | C+ |



Detection methods on Caltech-USA reasonable set

| Method | |
|---|---|
| VJ | 94.7% |
| Shapelet | 91.4% |
| PoseInv | 86.3% |
| LatSvm-V1 | 79.8% |
| ConvNet | 77.2% |
| FtrMine | 74.4% |
| HikSvm | 73.4% |
| HOG | 68.5% |
| MultiFtr | 68.3% |
| HogLbp | 67.8% |
| AFS+Geo | 66.8% |
| AFS | 65.4% |
| LatSvm-V2 | 63.3% |
| Pls | 62.1% |
| MLS | 61.0% |
| MultiFtr+CSS | 60.9% |
| FeatSynth | 60.2% |
| pAUCBoost | 59.7% |
| FPDW | 57.4% |
| ChnFtrs | 56.3% |
| CrossTalk | 53.9% |
| DBN−Isol | 53.1% |
| ACF | 51.4% |
| RandForest | 51.2% |
| MultiFtr+Motion | 50.9% |
| SquaresChnFtrs | 50.2% |
| Franken | 48.7% |
| MultiResC | 48.5% |
| Roerei | 48.4% |
| DBN−Mut | 48.2% |
| MF+Motion+2Ped | 46.4% |
| MOCO | 45.5% |
| MultiSDP | 45.4% |
| ACF-Caltech | 44.2% |
| MultiResC+2Ped | 43.4% |
| WordChannels | 42.3% |
| MT-DPM | 40.5% |
| JointDeep | 39.3% |
| SDN | 37.9% |
| MT-DPM+Context | 37.6% |
| ACF+SDt | 37.3% |
| SquaresChnFtrs | 34.8% |
| InformedHaar | 34.6% |
| Katamari-v1 | 22.5% |

INRIA training
Caltech-USA training
Other training

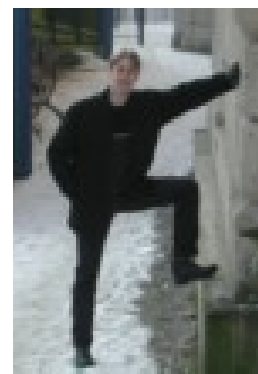log-average miss-rate (lower is better)

# Training data matters (you knew this already)

| Method | MR | Family | Features | Classifier | Context | Deep | Parts | M-Scales | More data | Feat. type | Training |
|---|---|---|---|---|---|---|---|---|---|---|---|
| VJ [9] | 94.73% | DF | ✓ | ✓ | | | | | | Haar | I |
| Shapelet [10] | 91.37% | - | | ✓ | | | | | | Gradients | I |
| PoseInv [11] | 86.32% | - | | | ✓ | | | | | HOG | I+ |
| LatSvm-V1 [12] | 79.78% | DPM | | | ✓ | | | | | HOG | P |
| ConvNet [13] | 77.20% | DN | | | | ✓ | | | | Pixels | I |
| FtrMine [14] | 74.42% | DF | ✓ | | | | | | | HOG+Color | I |
| HikSvm [15] | 73.39% | - | | ✓ | | | | | | HOG | I |
| HOG [1] | 68.46% | - | ✓ | ✓ | | | | | | HOG | I |
| MultiFtr [16] | 68.26% | DF | ✓ | ✓ | | | | | | HOG+Haar | I |
| HogLbp [17] | 67.77% | - | | ✓ | | | | | | HOG+LBP | I |
| AFS+Geo [18] | 66.76% | - | | | ✓ | | | | | Custom | I |
| AFS [18] | 65.38% | - | | | | | | | | Custom | I |
| LatSvm-V2 [19] | 63.26% | DPM | ✓ | | | | ✓ | | | HOG | I |
| Pls [20] | 62.10% | - | ✓ | ✓ | | | | | | Custom | I |
| MLS [21] | 61.03% | DF | ✓ | | | | | | | HOG | I |
| MultiFtr+CSS [22] | 60.89% | DF | ✓ | | | | | | | Many | T |
| FeatSynth [23] | 60.16% | - | ✓ | ✓ | | | | | | Custom | I |
| pAUCBoost [24] | 59.66% | DF | ✓ | ✓ | | | | | | HOG+COV | I |
| FPDW [25] | 57.40% | DF | | | | | | | | HOG+LUV | I |
| ChnFtrs [26] | 56.34% | DF | ✓ | ✓ | | | | | | HOG+LUV | I |
| CrossTalk [27] | 53.88% | DF | | | ✓ | | | | | HOG+LUV | I |
| DBN-Isol [28] | 53.14% | DN | | | | | ✓ | | | HOG | I |
| ACF [29] | 51.36% | DF | ✓ | | | | | | | HOG+LUV | I |
| RandForest [30] | 51.17% | DF | | ✓ | | | | | | HOG+LBP | I&C |
| MultiFtr+Motion [22] | 50.88% | DF | ✓ | | | | | | ✓ | Many+Flow | T |
| SquaresChnFtrs [31] | 50.17% | DF | ✓ | | | | | | | HOG+LUV | I |
| Franken [32] | 48.68% | DF | | ✓ | | | | | | HOG+LUV | I |
| MultiResC [33] | 48.45% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C |
| Roerei [31] | 48.35% | DF | ✓ | | | | | ✓ | | HOG+LUV | I |
| DBN-Mut [34] | 48.22% | DN | | | ✓ | | ✓ | | | HOG | C |
| MF+Motion+2Ped [35] | 46.44% | DF | | | ✓ | | | | ✓ | Many+Flow | I+ |
| MOCO [36] | 45.53% | - | ✓ | | ✓ | | | | | HOG+LBP | C |
| MultiSDP [37] | 45.39% | DN | ✓ | | ✓ | ✓ | | | | HOG+CSS | C |
| ACF-Caltech [29] | 44.22% | DF | ✓ | | | | | | | HOG+LUV | C |
| MultiResC+2Ped [35] | 43.42% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C+ |
| WordChannels [38] | 42.30% | DF | ✓ | | | | | | | Many | C |
| MT-DPM [39] | 40.54% | DPM | | | | | ✓ | ✓ | | HOG | C |
| JointDeep [40] | 39.32% | DN | | | | ✓ | | | | Color+Gradient | C |
| SDN [41] | 37.87% | DN | | | | ✓ | ✓ | | | Pixels | C |
| MT-DPM+Context [39] | 37.64% | DPM | | | ✓ | | ✓ | ✓ | | HOG | C+ |
| ACF+SDt [42] | 37.34% | DF | ✓ | | | | | | ✓ | ACF+Flow | C+ |
| SquaresChnFtrs [31] | 34.81% | DF | ✓ | | | | | | | HOG+LUV | C |
| InformedHaar [43] | 34.60% | DF | ✓ | | | | | | | HOG+LUV | C |
| Katamari-v1 | 22.49% | DF | ✓ | | ✓ | | | | ✓ | HOG+Flow | C+ |

## Detection methods on Caltech-USA reasonable set

VJ 94.7%
Shapelet 91.4%
PoseInv 86.3%
LatSvm-V1 79.8%
ConvNet 77.2%
FtrMine 74.4%
HikSvm 73.4%
HOG 68.5%
MultiFtr 68.3%
HogLbp 67.8%
AFS+Geo 66.8%
AFS 65.4%
LatSvm-V2 63.3%
Pls 62.1%
MLS 61.0%
MultiFtr+CSS 60.9%
FeatSynth 60.2%
pAUCBoost 59.7%
FPDW 57.4%
ChnFtrs 56.3%
CrossTalk 53.9%
DBN−Isol 53.1%
ACF 51.4%
RandForest 51.2%
MultiFtr+Motion 50.9%
SquaresChnFtrs 50.2%
Franken 48.7%
MultiResC 48.5%
Roerei 48.4%
DBN−Mut 48.2%
MF+Motion+2Ped 46.4%
MOCO 45.5%
MultiSDP 45.4%
ACF-Caltech 44.2%
MultiResC+2Ped 43.4%
WordChannels 42.3%
MT-DPM 40.5%
JointDeep 39.3%
SDN 37.9%
MT-DPM+Context 37.6%
ACF+SDt 37.3%
SquaresChnFtrs 34.8%
InformedHaar 34.6%
Katamari-v1 22.5%

Legend:
- INRIA training
- Caltech-USA training
- Other training

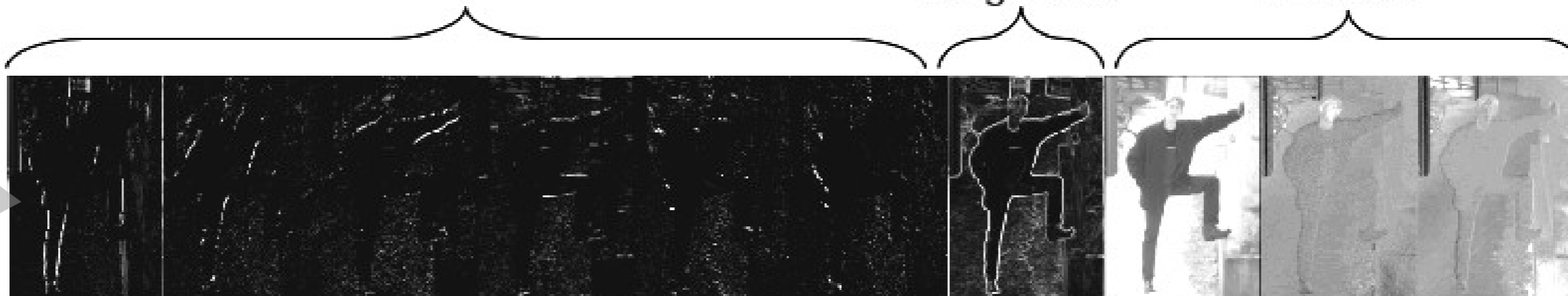log-average miss-rate (lower is better)

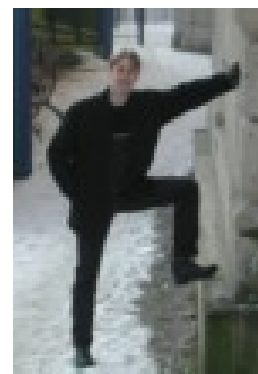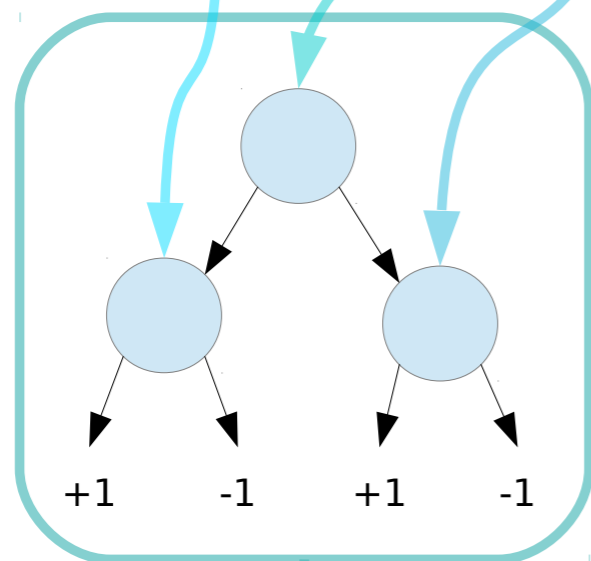# Strong detection with (shallow) boosted decision trees



6 Orientation bins      Gradient magnitude      LUV color channels

# Strong detection with (shallow) boosted decision trees
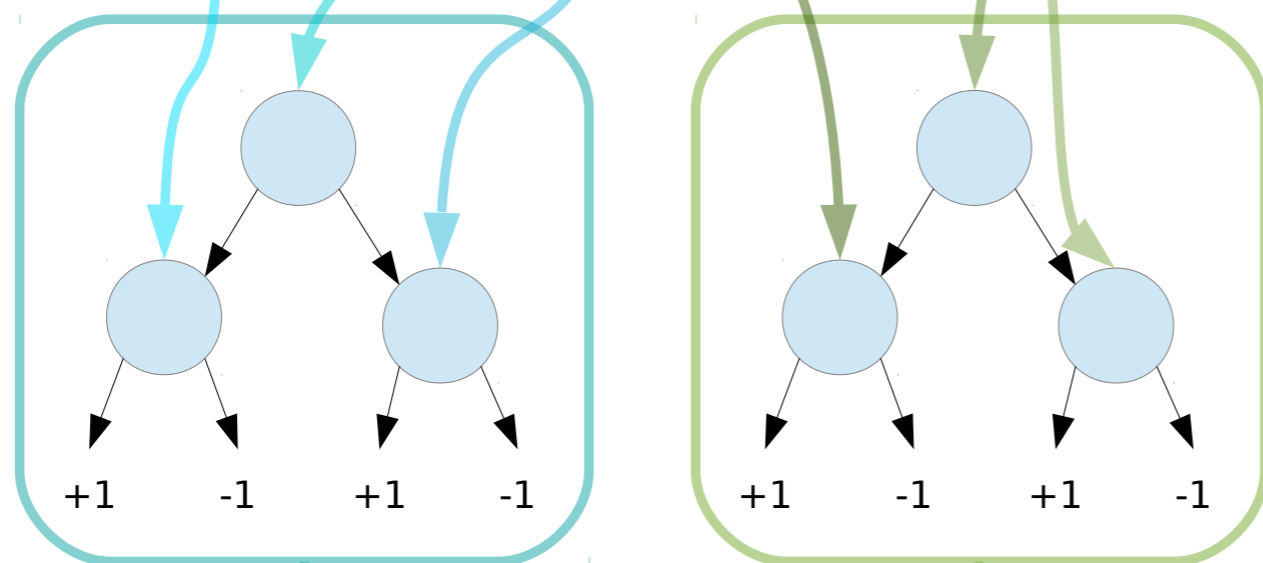


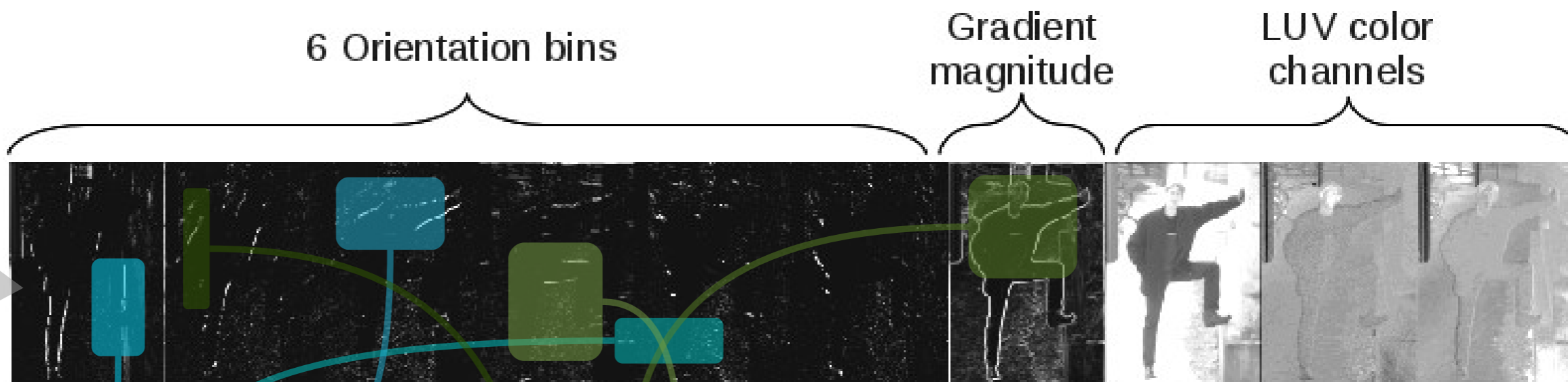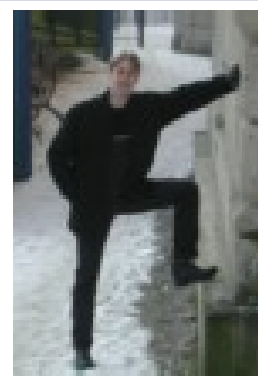6 Orientation bins     Gradient magnitude     LUV color channels

+1   -1   +1   -1

$$score = \; w_1 \cdot h_1 +$$

# Strong detection with (shallow) boosted decision trees



6 Orientation bins     Gradient magnitude     LUV color channels

+1   -1   +1   -1     +1   -1   +1   -1

$$score = \; w_1 \cdot h_1 + \qquad\qquad w_2 \cdot h_2 +$$

# Strong detection with (shallow) boosted decision trees



6 Orientation bins

Gradient magnitude

LUV color channels

+1  -1  +1  -1

+1  -1  +1  -1

• • •

+1  -1  +1  -1

$$score = w_1 \cdot h_1 + \quad w_2 \cdot h_2 + \quad \cdots \quad + w_N \cdot h_N$$

[ChnFtrs, Dollar et al. 2009; SquaresChnFtrs, Benenson et al. 2013]

# Strong detection with (shallow) boosted decision trees



6 Orientation bins

Gradient magnitude

LUV color channels

+1  -1  +1  -1

+1  -1  +1  -1

• • •

+1  -1  +1  -1

**"Viola&Jones meets Dalal&Triggs"** (2001 & 2005)

[ChnFtrs, Dollar et al. 2009; SquaresChnFtrs, Benenson et al. 2013]

# Only pedestrians ?

**Video at http://goo.gl/14cz07**

[Mathias et al. IJCNN 2013] [Mathias et al. ECCV 2014]

**Video at http://goo.gl/Evayrz**

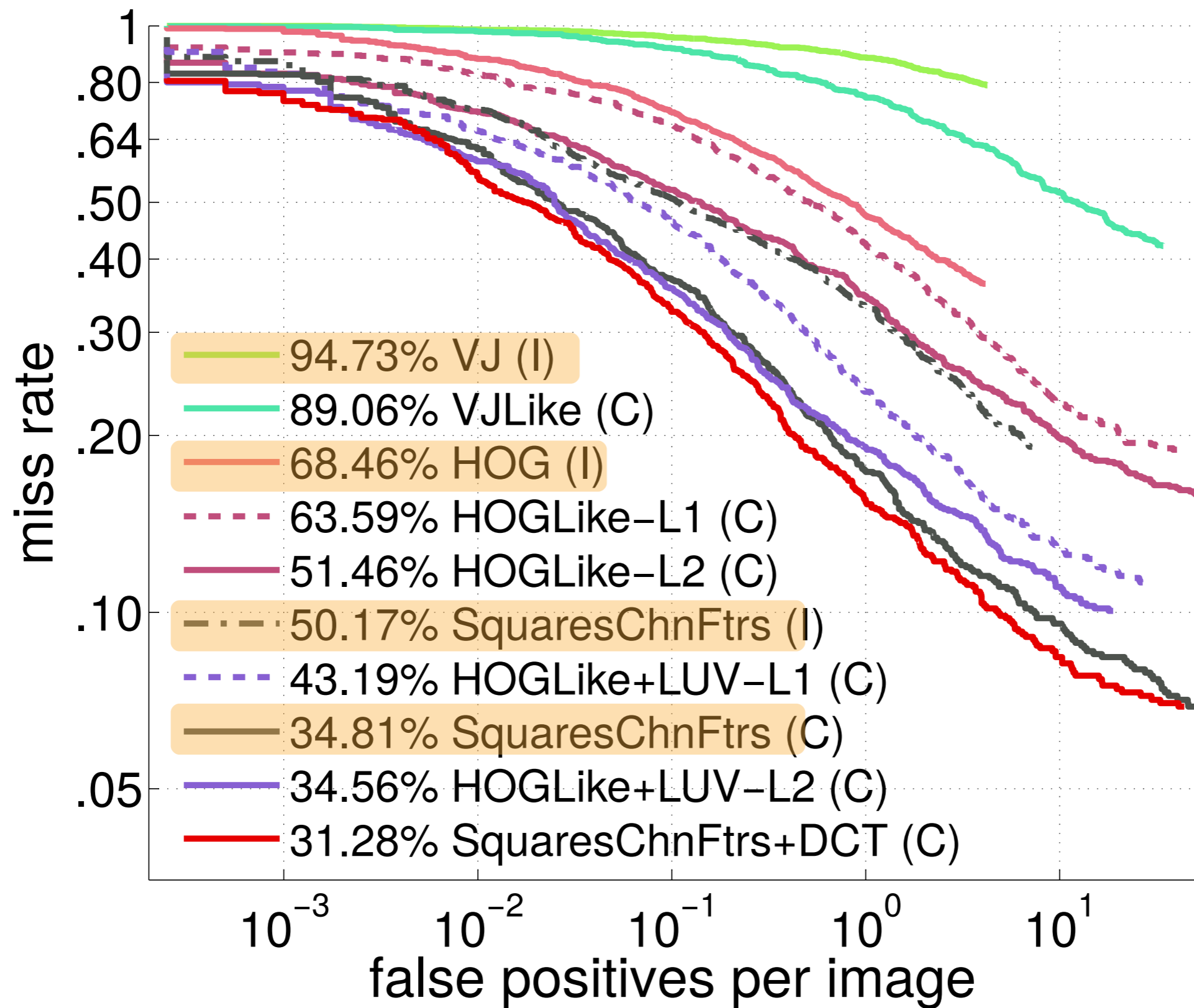[Mathias et al. IJCNN 2013] [Mathias et al. ECCV 2014]

# What is driving the quality progress ?

- ~~solution family (DPM, deep networks, decision forests)~~

- ~~better classifiers~~

- ~~deformable parts~~

- ~~multi-scale models~~

- ~~deep architectures~~

- ~~training data~~

- **additional (test time) data**
  ⇨ using more frames (flow or stereo) helps (you knew this already)
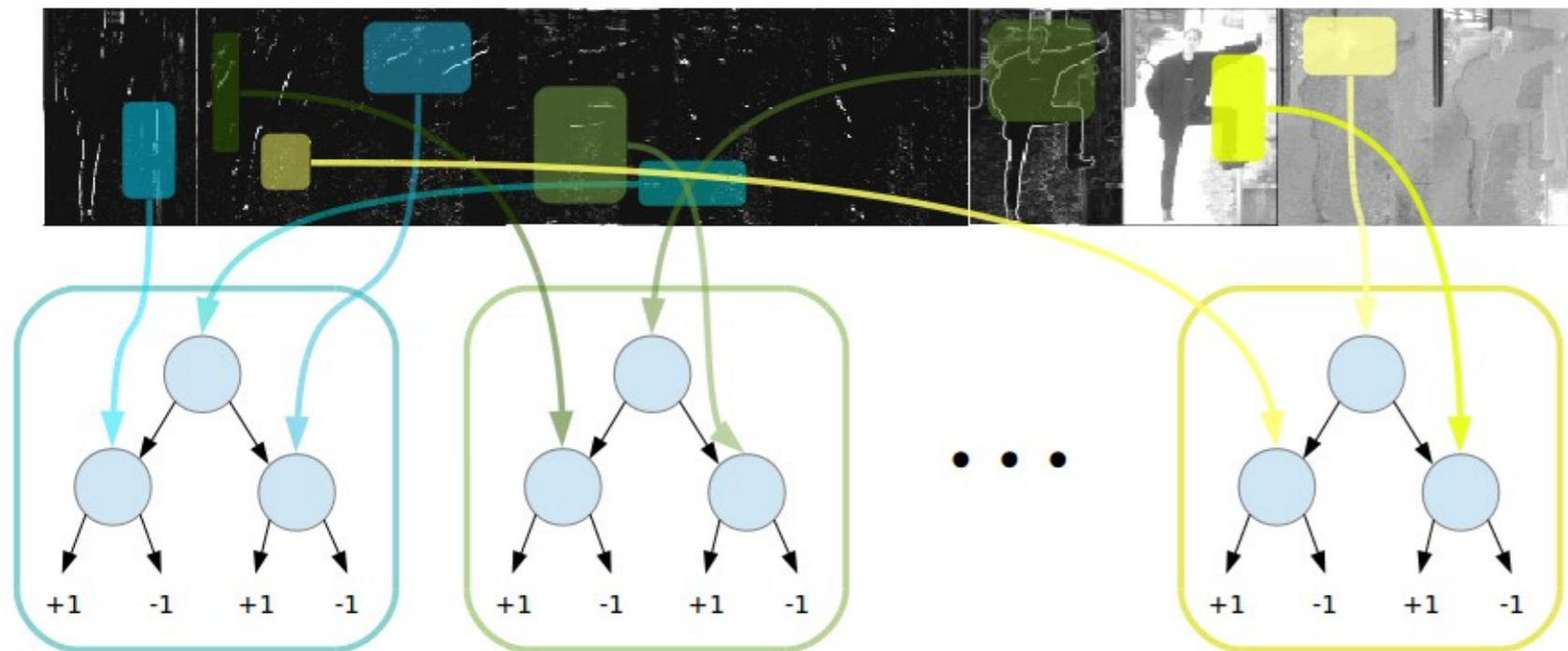
- exploiting context

- better features

# What is driving the quality progress ?

- ~~solution family (DPM, deep networks, decision forests)~~

- ~~better classifiers~~

- ~~deformable parts~~

- ~~multi-scale models~~

- ~~deep architectures~~

- ~~training data~~

- ~~additional (test time) data~~

- **exploiting context**

- better features

# Using context helps (expect ~5 pp improvement)



Relative improvement on Caltech-USA reasonable set

# What is driving the quality progress ?

- ~~solution family (DPM, deep networks, decision forests)~~

- ~~better classifiers~~

- ~~deformable parts~~

- ~~multi-scale models~~

- ~~deep architectures~~

- ~~training data~~

- ~~additional (test time) data~~

- ~~exploiting context~~

- **better features**

# Experiments

(some of them)

# Features alone can explain 10 years of progress



Legend (miss rate, false positives per image):
- 94.73% VJ (I)
- 89.06% VJLike (C)
- 68.46% HOG (I)
- 63.59% HOGLike–L1 (C)
- 51.46% HOGLike–L2 (C)
- 50.17% SquaresChnFtrs (I)
- 43.19% HOGLike+LUV–L1 (C)
- 34.81% SquaresChnFtrs (C)
- 34.56% HOGLike+LUV–L2 (C)
- 31.28% SquaresChnFtrs+DCT (C)

Axes: miss rate vs. false positives per image

# What is driving the quality progress ?

- solution family (DPM, deep networks, decision forests)

- better classifiers

- deformable parts

- multi-scale models

- deep architectures

- training data

- additional (test time) data

- exploiting context

- better features

# Strong features/Flow/Context are very complementary

| Method | Results | Improvement | Expected improvement |
|---|---|---|---|
| SquaresChnFtrs | 34.81% | - | - |

Results in MR (lower is better). Improvement in MR percent points.



[Benenson et al. 2013]

# Strong features/Flow/Context are very complementary

| Method | Results | Improvement | Expected improvement |
|---|---|---|---|
| SquaresChnFtrs | 34.81% | - | - |
| +Better features (DCT) | 31.28% | 3.53 | - |

Results in MR (lower is better). Improvement in MR percent points.



**\***
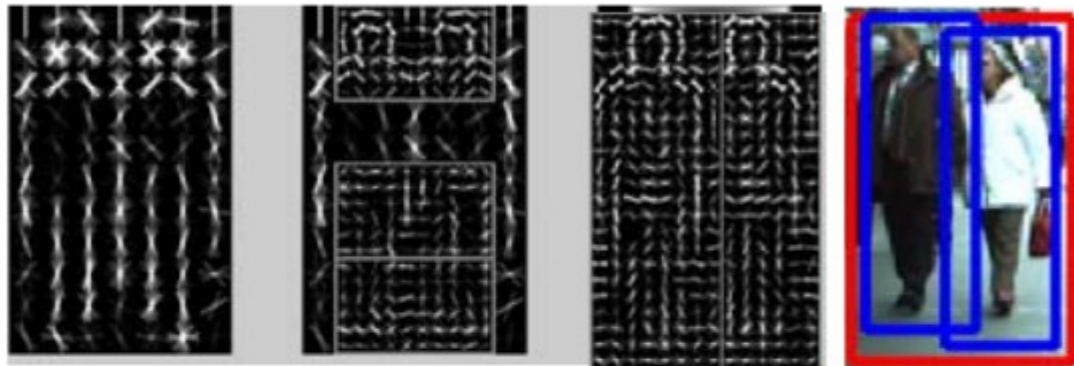
[DCT: Nam et al. ArXiv 2014]

# Strong features/Flow/Context are very complementary

| Method | Results | Improvement | Expected improvement |
|---|---|---|---|
| SquaresChnFtrs | 34.81% | - | - |
| +Better features (DCT) | 31.28% | 3.53 | - |
| +Flow (SDt) | 30.34% | 4.47 | - |

Results in MR (lower is better). Improvement in MR percent points.



Stabilized, m=1    Stabilized, m=8

t=0    t=8, stabilized

[DCT: Nam et al. ArXiv 2014]
[SDt: Park et al. CVPR 2013]

# Strong features/Flow/Context are very complementary

| Method | Results | Improvement | Expected improvement |
|---|---|---|---|
| SquaresChnFtrs | 34.81% | - | - |
| +Better features (DCT) | 31.28% | 3.53 | - |
| +Flow (SDt) | 30.34% | 4.47 | - |
| +Context (2Ped) | 29.42% | 5.39 | - |

Results in MR (lower is better). Improvement in MR percent points.



Aspect Ratio 2: 12x7

[DCT: Nam et al. ArXiv 2014]
[SDt: Park et al. CVPR 2013]
[2Ped: Ouyang & Wang CVPR 2013]

# Strong features/Flow/Context are very complementary

| Method | Results | Improvement | Expected improvement |
|---|---|---|---|
| SquaresChnFtrs | 34.81% | - | - |
| +Better features (DCT) | 31.28% | 3.53 | - |
| +Flow (SDt) | 30.34% | 4.47 | - |
| +Context (2Ped) | 29.42% | 5.39 | - |

Results in MR (lower is better). Improvement in MR percent points.

[DCT: Nam et al. ArXiv 2014]
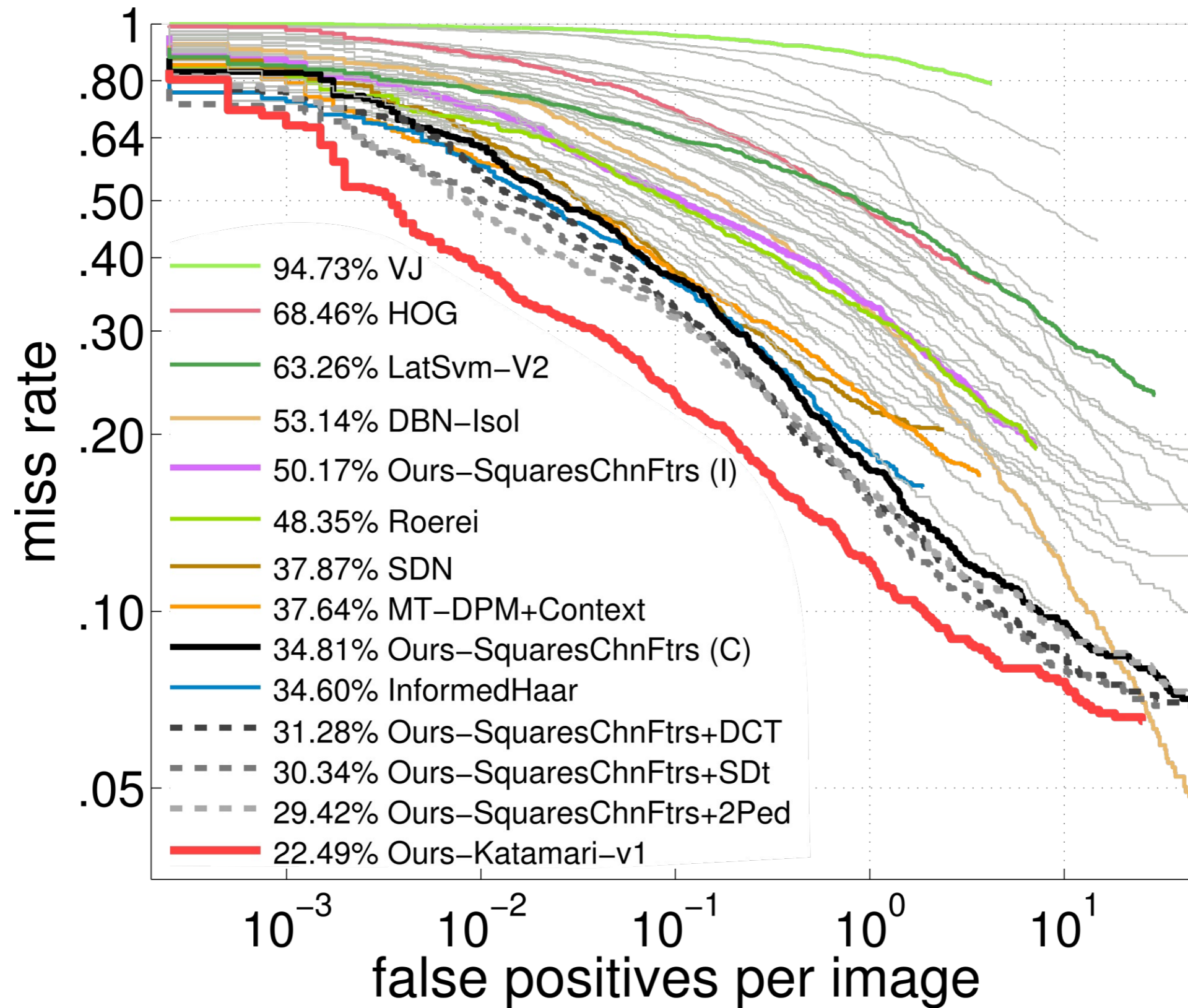[SDt: Park et al. CVPR 2013]
[2Ped: Ouyang & Wang CVPR 2013]

# Strong features/Flow/Context are very complementary

| Method | Results | Improvement | Expected improvement |
|---|---|---|---|
| SquaresChnFtrs | 34.81% | - | - |
| +Better features (DCT) | 31.28% | 3.53 | - |
| +Flow (SDt) | 30.34% | 4.47 | - |
| +Context (2Ped) | 29.42% | 5.39 | - |
| +DCT+2Ped | 27.40% | 7.41 | 8.92 |
| +SDt+2Ped | 26.68% | 8.13 | 9.86 |
| +DCT+SDt | 25.24% | 9.57 | 8.00 |
| All-in-one (Katamari) | 22.49% | 12.32 | 13.39 |

Results in MR (lower is better). Improvement in MR percent points.

Surprise 2: no diminishing return observed (yet).

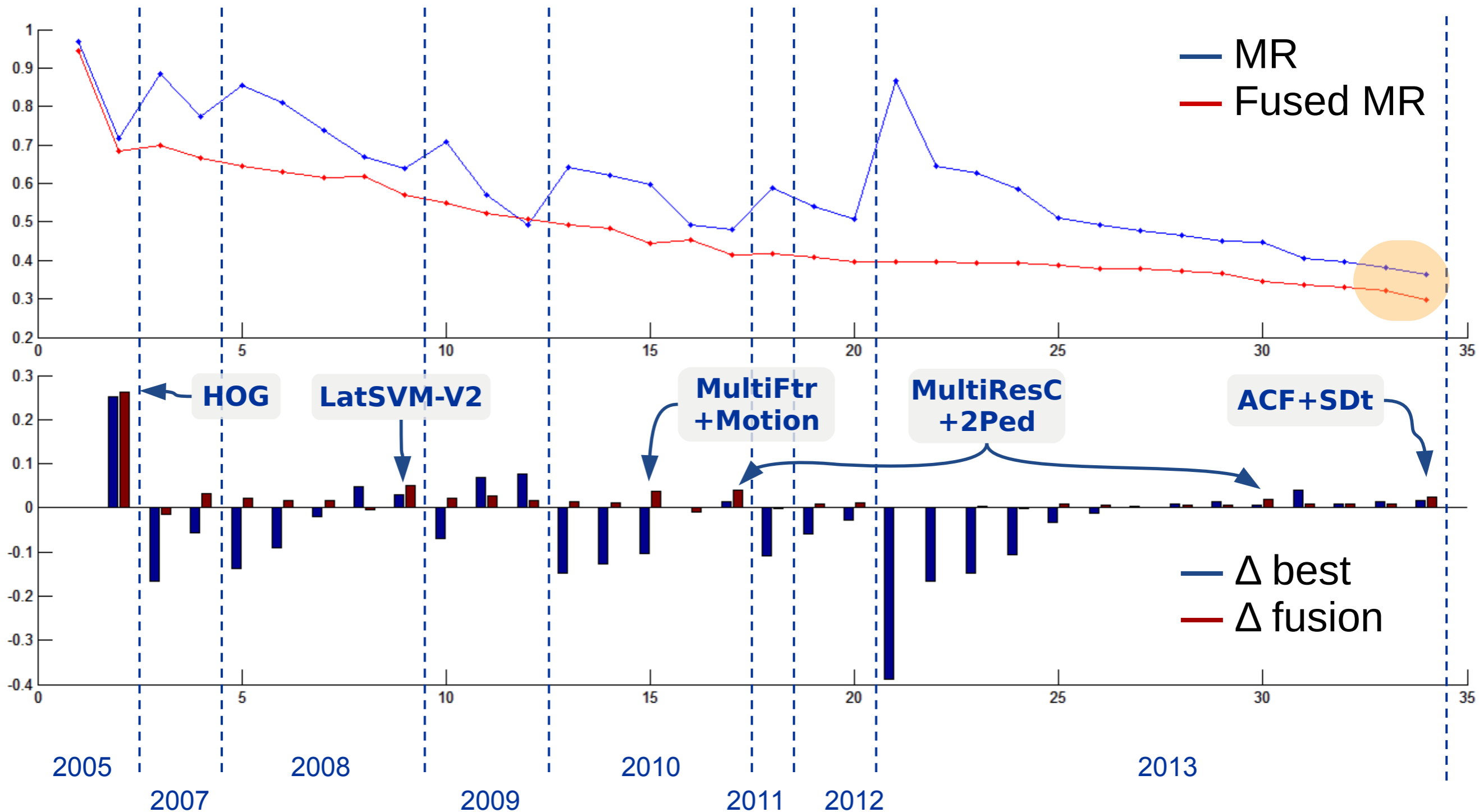# Strong features/Flow/Context are very complementary



Legend (miss rate):
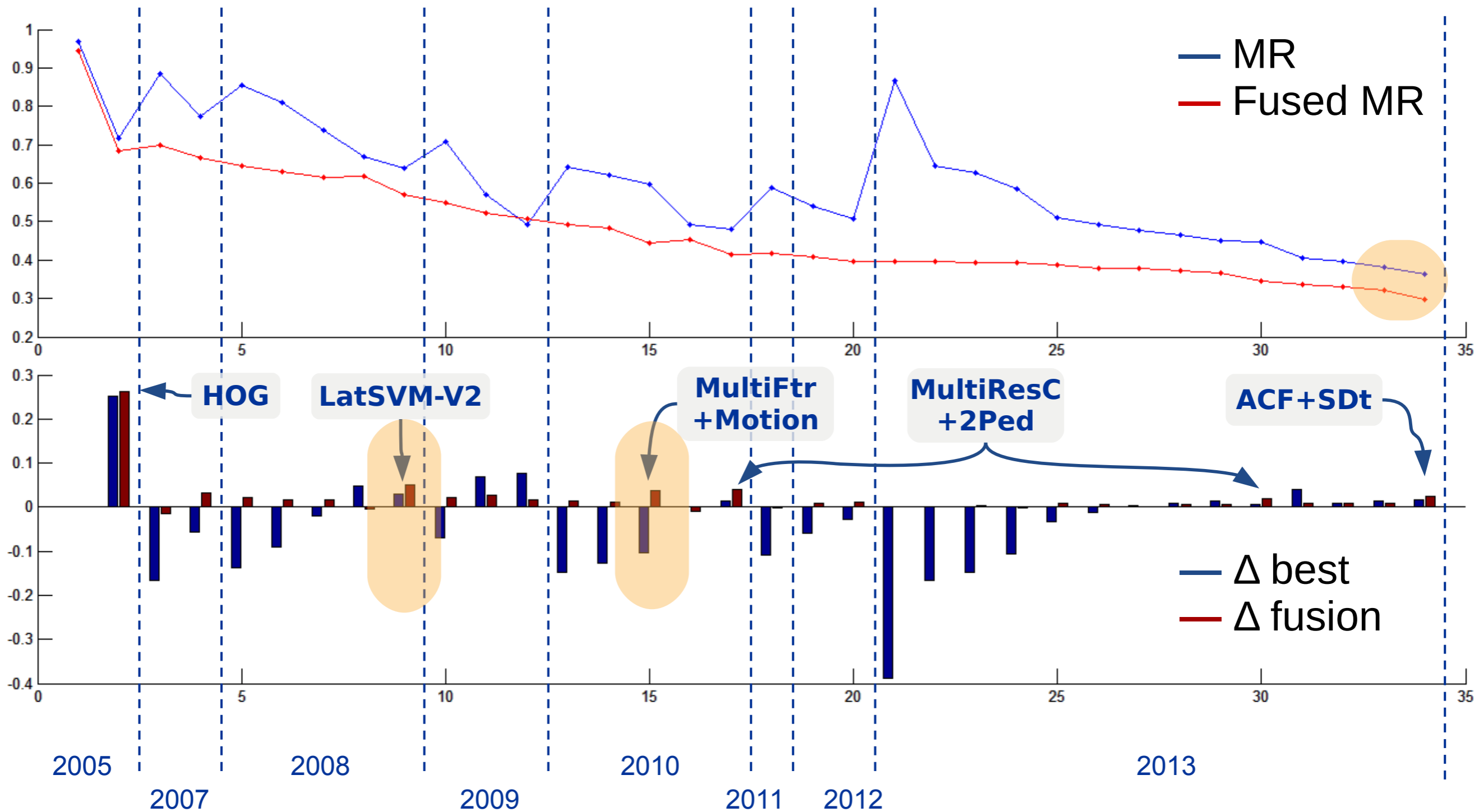- 94.73% VJ
- 68.46% HOG
- 63.26% LatSvm−V2
- 53.14% DBN−Isol
- 50.17% Ours−SquaresChnFtrs (I)
- 48.35% Roerei
- 37.87% SDN
- 37.64% MT−DPM+Context
- 34.81% Ours−SquaresChnFtrs (C)
- 34.60% InformedHaar
- 31.28% Ours−SquaresChnFtrs+DCT
- 30.34% Ours−SquaresChnFtrs+SDt
- 29.42% Ours−SquaresChnFtrs+2Ped
- 22.49% Ours−Katamari−v1

miss rate (y-axis)

false positives per image (x-axis)

# Merging all methods over time



Slide from [Xu et al. BMVC 2014]

# Merging all methods over time



HOG

LatSVM-V2

MultiFtr +Motion

MultiResC +2Ped

ACF+SDt

MR
Fused MR

Δ best
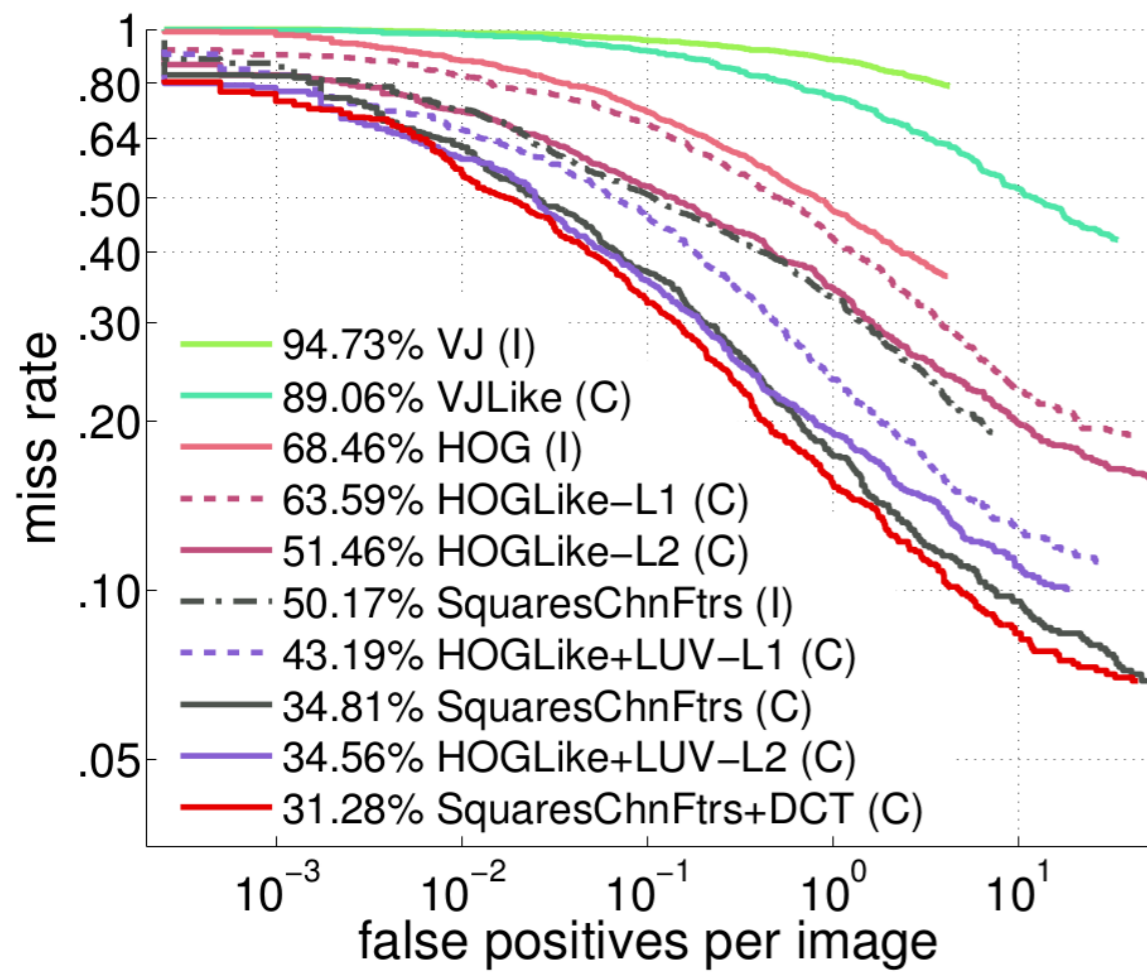Δ fusion

2005
2007
2008
2009
2010
2011
2012
2013

Slide from [Xu et al. BMVC 2014]
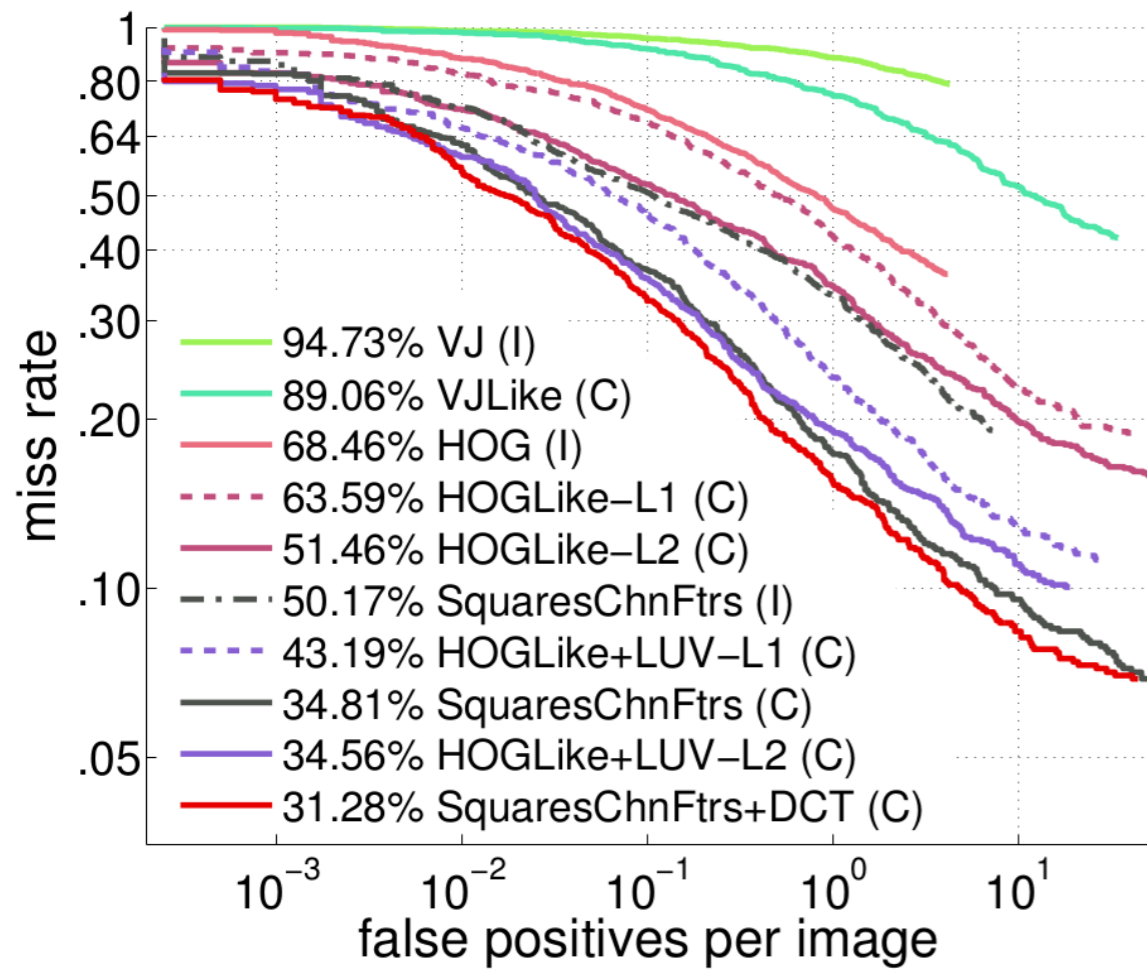
# Merging all methods over time



Slide from [Xu et al. BMVC 2014]
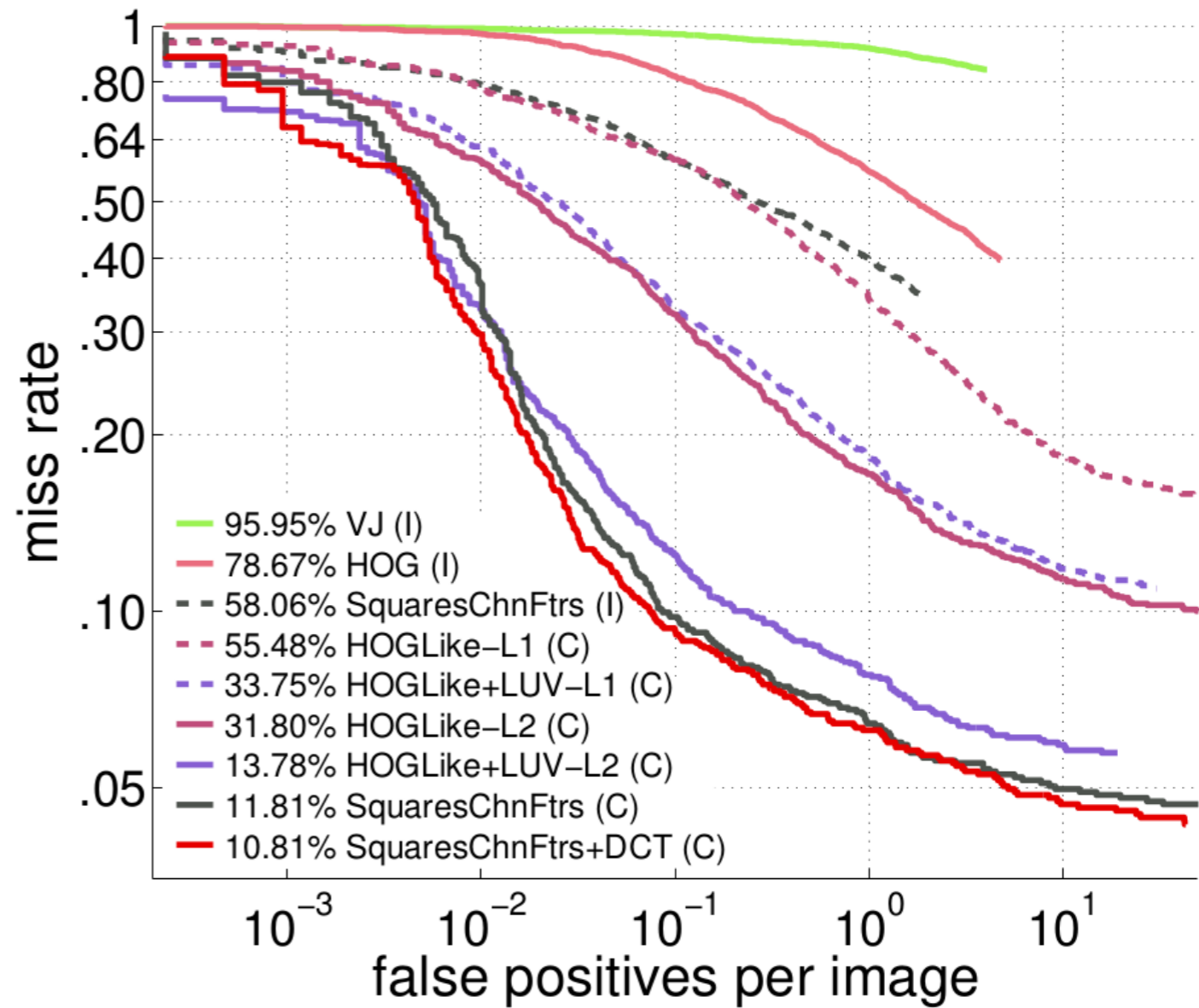
# Surprise 3: Model capacity has not saturated



Caltech-USA
test set

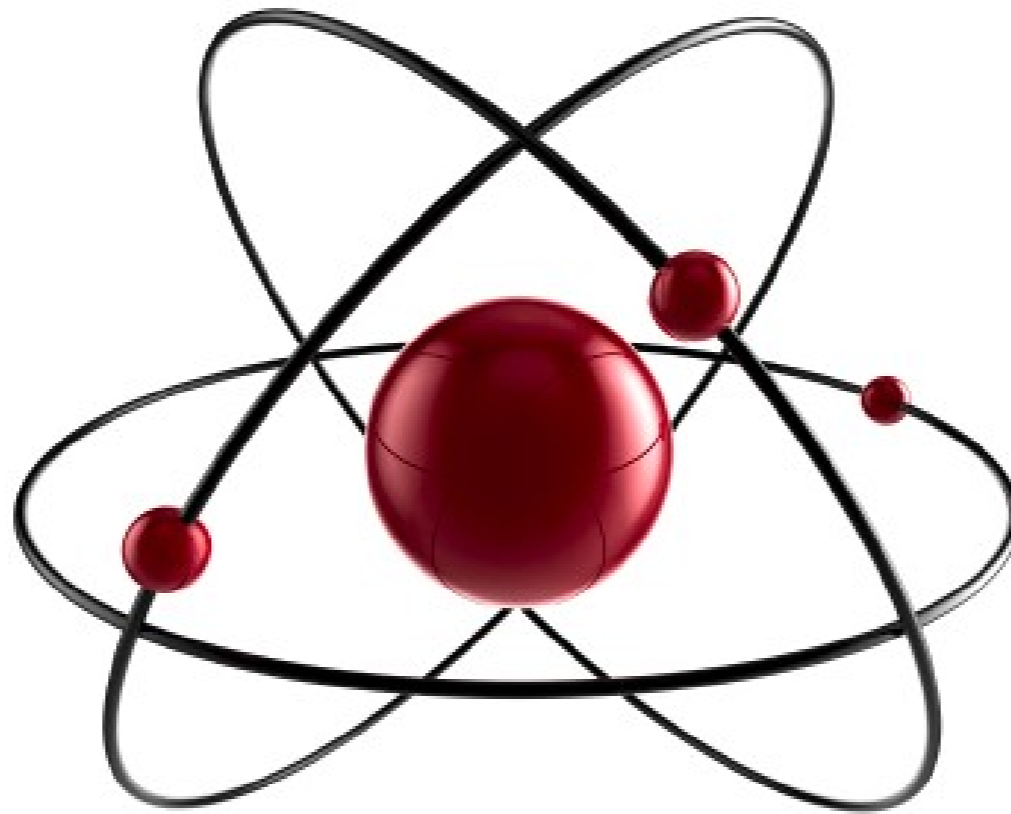# Surprise 3: Model capacity has not saturated



Caltech-USA
test set

Caltech-USA
**training set**

# What have we learned ?

- "Sooner or later, everything old is new again." - Stephen King
  Decade-old ideas still rule detection quality.

- Switching training data is not comparing apples-to-apples.

- Flow, context, and strong features are very complementary (still).

- All other aspects have yet to make a "definitive statement".

- Features alone can explain a decade of detection quality progress.

- There is room for further improvement
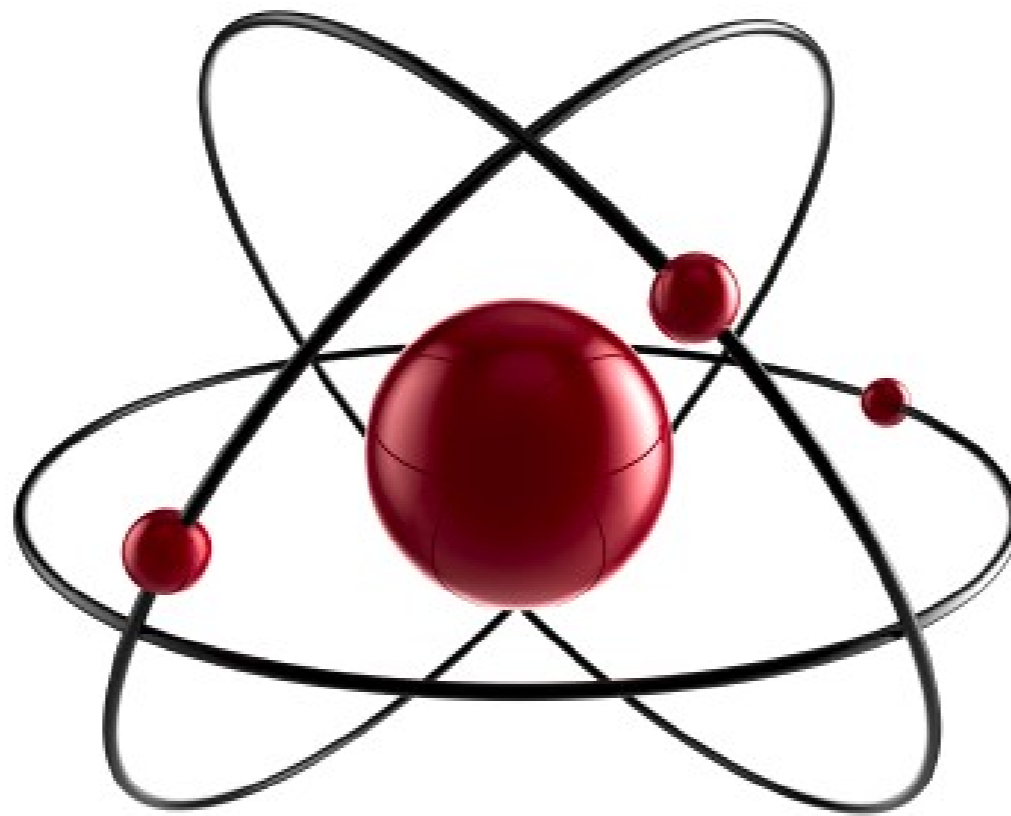  by increasing model capacity (and better features).

Message of the day:

One (simple and effective) core

+

3 add-ons

Message of the day:

"Viola&Jones meets Dalal&Triggs"

+

Better features + Context + Flow

# How to further improve quality ?

- **Stronger use of additional data**
  (scene flow on KITTI ?)

- **Better context**
  (exploiting scene geometry)

- **Further developing deep architectures**
  (end-to-end fine tuning)

- **Most importantly: understanding**
  what makes good features good?