

# Class-Specific Intra-Frame Tracking for Fast Detectors

Fatih Porikli\*

Oncel Tuzel†

\*Mitsubishi Electric Research Laboratories, USA, †Rutgers University, USA

**Don't do multi-scale, multi-rotation, multi-warp, affine, etc. object search. Too expensive! Do Intra-frame tracking. We do it with regression. TD rate 5% → 96.7%.**

$$M = \begin{pmatrix} A & b \\ 0 & 1 \end{pmatrix} \quad M_t = M_{t-1} \cdot \Delta M_t \quad \Delta M_t = f(\mathbf{o}_t(M_{t-1}^{-1}))$$

$$\exp(\mathbf{m}) = M$$

$$\rho(M_1, M_2) = \|\log(M_1^{-1}M_2)\| \approx \|\mathbf{m}_2 - \mathbf{m}_1\|$$

$$f(\mathbf{o}) = \exp(g(\mathbf{o})) \quad g(\mathbf{o}) = \mathbf{o}^T \Omega$$

$$J_g = \sum_{i=1}^n \rho^2 [f(\mathbf{o}_0^i), \Delta M_i] \quad \Delta M_t = f(\mathbf{o}_t(M_{t-1}^{-1}))$$

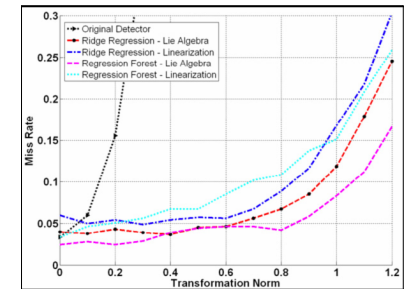
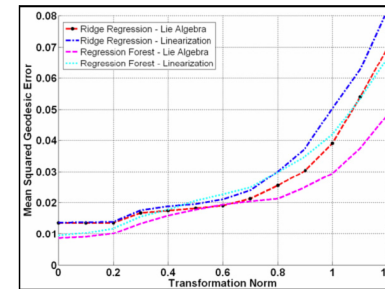
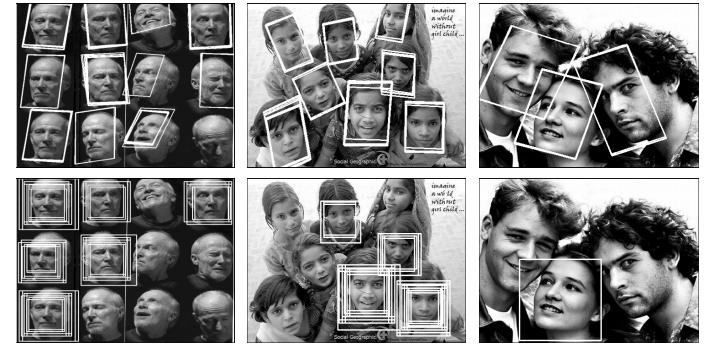
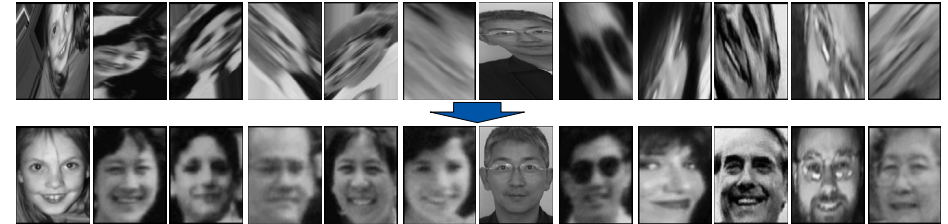
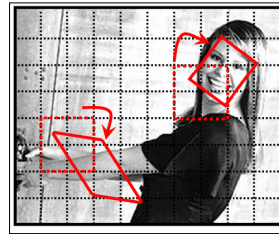
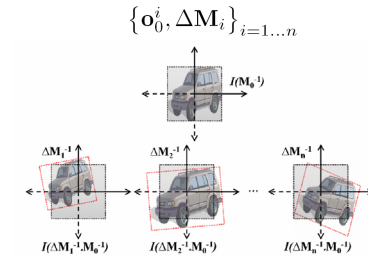
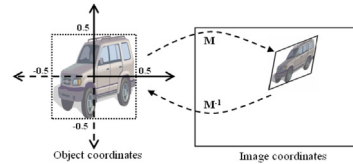
$$J_a = \sum_{i=1}^n \|\log(f(\mathbf{o}_0^i)) - \log(\Delta M_i)\|^2$$

$$X = \begin{pmatrix} [\mathbf{o}_0^1]^T \\ \vdots \\ [\mathbf{o}_0^n]^T \end{pmatrix} \quad Y = \begin{pmatrix} [\log(\Delta M_1)]^T \\ \vdots \\ [\log(\Delta M_n)]^T \end{pmatrix}$$

$$J_a = \text{tr}[(X\Omega - Y)^T(X\Omega - Y)]$$

**Input:** Location of target at time  $t-1$  is  $M_{t-1}$  and the current observation is  $I_t$ , maximum iteration number is  $K$ .

- $k = 1$  and  $M_t = M_{t-1}$
- Repeat
  - $\Delta M_t = f(\mathbf{o}_t(M_{t-1}^{-1}))$
  - $M_t = M_{t-1} \cdot \Delta M_t$
  - $k = k + 1$
- Until  $\Delta M_t = I$  or  $k = K$



$$J_r = \text{tr}[(X\Omega - Y)^T(X\Omega - Y)] + \lambda \|\Omega\|^2$$

$$\Omega = (X^T X + \lambda I)^{-1} X^T Y$$

•803 face images from CMU, MIT and MERL datasets.

•The dataset is divided into 503 images for training and 300 for testing.

•The training set consists of 25150 samples which are generated by applying 50 transformations having a random norm between 0 to 1.0 to each face image



fatih@merl.com

•For an image of size 320x240, the VJ detector evaluates 58367 locations for translation and scale search, whereas the proposed method evaluates face detector at only 642 locations searched on the affine space.

•Tracker iterations are performed which requires 12840 warping operations.

•Since the warps at each iteration can be performed in parallel we implemented the warping in GPU using NVIDIA GeForce 8800 GTX graphics card and CUDA SDK.

•The search for an image of size 320x240 takes 0.85 and 2.4 seconds with the linear and the regression forest models respectively. On average, Lie algebra based parametrization have 50% less miss rate for regression forest and 24% less for ridge regression, compared to linearization.