

# Gaussian Process II: Learning hyper parameters

## Choosing kernel parameters

$$k(t_i, t_j) = \sigma_f^2 \exp \left\{ -\frac{1}{2\ell^2} (t_i - t_j)^2 \right\}$$

- ▶ Choose some *hyperparameters*:  $\sigma_f = 14$ ,  $\ell = 50$

$$t = \begin{bmatrix} 0700 \\ 0800 \\ 1029 \end{bmatrix} \quad K(t, t) = \{k(t_i, t_j)\}_{i,j} = \begin{bmatrix} 196 & 26.5 & 00.0 \\ 26.5 & 196 & 0.01 \\ 00.0 & 0.01 & 196 \end{bmatrix}$$

- ▶ Choose some *hyperparameters*:  $\sigma_f = 7$ ,  $\ell = 100$

$$t = \begin{bmatrix} 0700 \\ 0800 \\ 1029 \end{bmatrix} \quad K(t, t) = \{k(t_i, t_j)\}_{i,j} = \begin{bmatrix} 49.0 & 29.7 & 0.2 \\ 29.7 & 49.0 & 03.6 \\ 00.2 & 03.6 & 49.0 \end{bmatrix}$$

- ▶ Choose some *hyperparameters*:  $\sigma_f = 7$ ,  $\ell = 500$

$$t = \begin{bmatrix} 0700 \\ 0800 \\ 1029 \end{bmatrix} \quad K(t, t) = \{k(t_i, t_j)\}_{i,j} = \begin{bmatrix} 49.0 & 48.0 & 39.5 \\ 48.0 & 49.0 & 44.1 \\ 39.5 & 44.1 & 49.0 \end{bmatrix}$$

## Sampling from a Covariance matrix

Scalar case

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

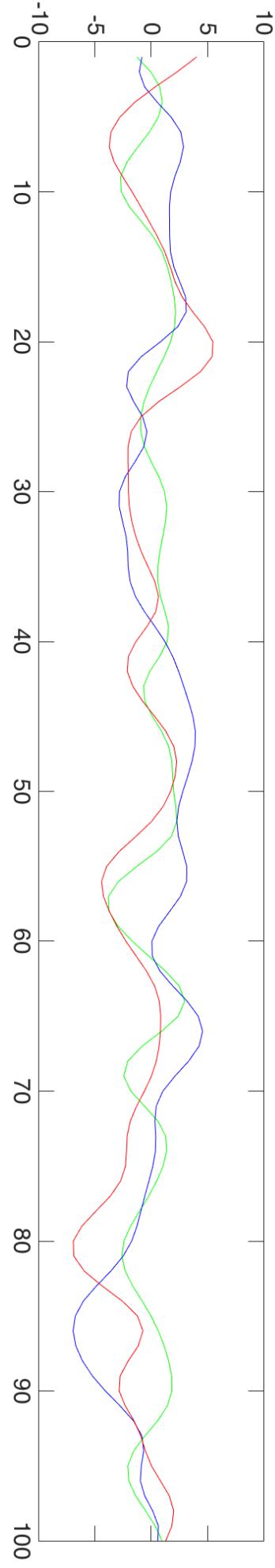
$$x \sim \mu + \sigma(\mathcal{N}(0, 1))$$

Vector case in GP setting

$$f_* \sim \mu + B(\mathcal{N}(0, I))$$

$$BB^T = \Sigma_*$$

Cholesky decomposition accomplishes this factorization  
where are  $B$  and  $B^T$  upper and lower triangular matrices



```
function [ Km ] = GaussianF( numPts,s,l )
%Priors for GP; numPts=100, s=2, l=3

Km = zeros(numPts,numPts);
for i = 1: numPts
    for j = 1: numPts
        Km(i,j) = (s^2)*exp(-(0.5/l^2)*(i-j)^2);
    end
end
end

F=chol(Km)'*randn(100,1);
```

If  $f$  and  $y$  are jointly Gaussian:

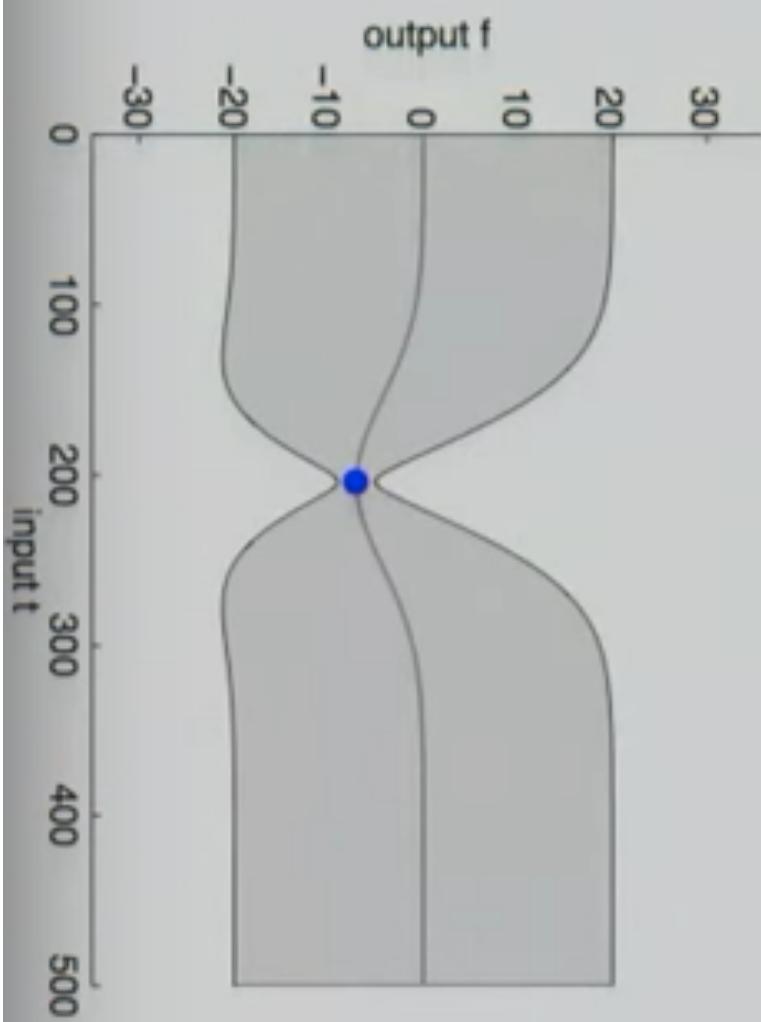
$$\begin{bmatrix} f \\ y \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} m_f \\ m_y \end{bmatrix}, \begin{bmatrix} K_{ff} & K_{fy} \\ K_{fy}^T & K_{yy} \end{bmatrix} \right)$$

Then:

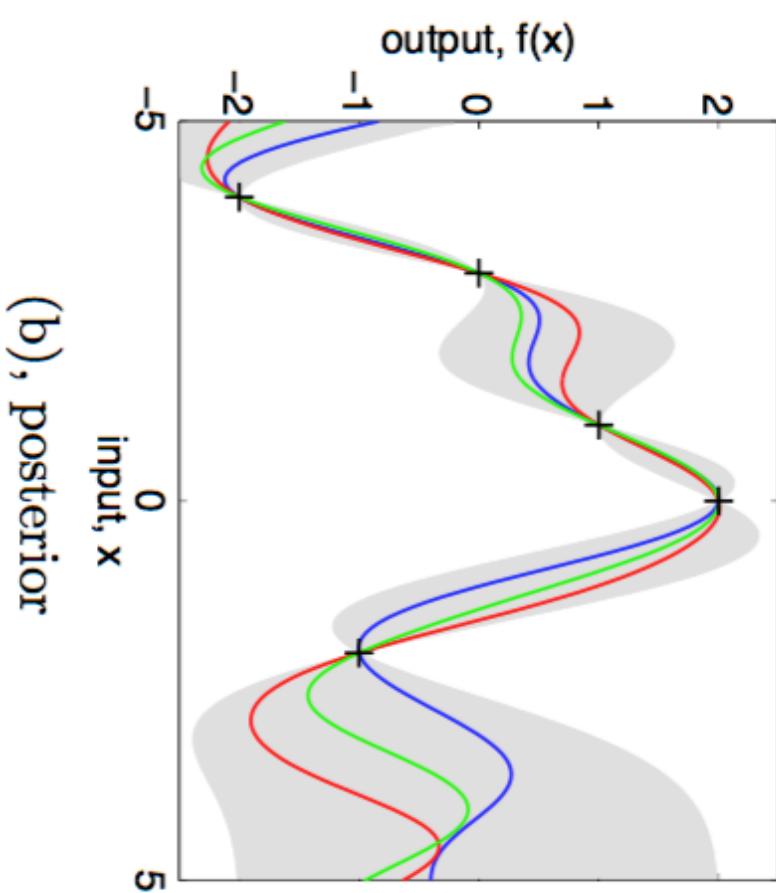
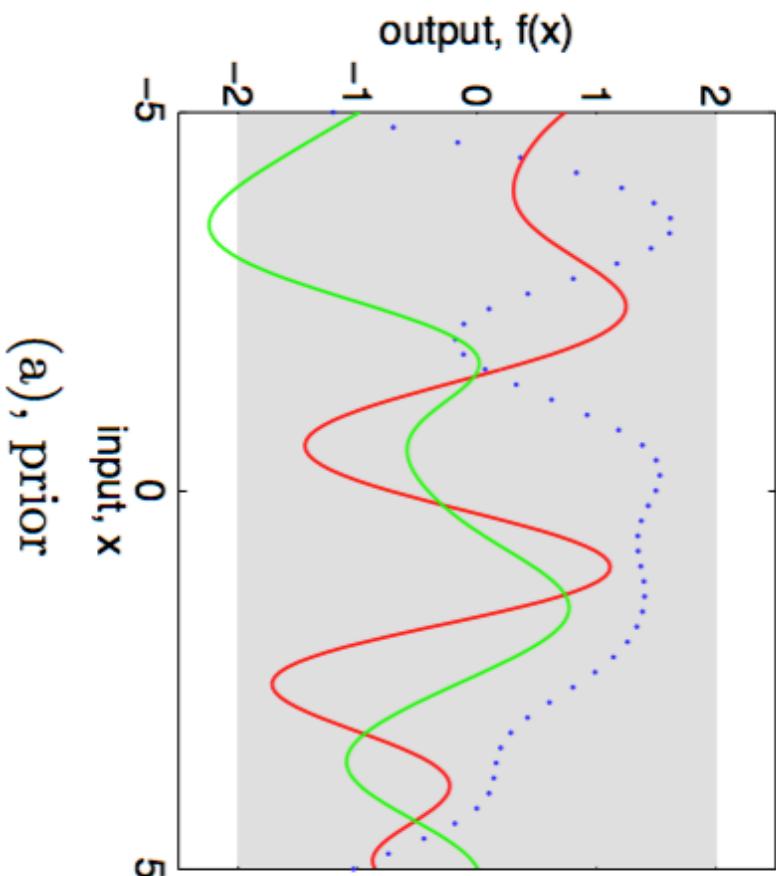
$$f|y \sim \mathcal{N} (K_{fy} K_{yy}^{-1} (y - m_y) + m_f, K_{ff} - K_{fy} K_{yy}^{-1} K_{fy}^T)$$

► Use conditioning to update the posterior:

$$f|y(204) \sim \mathcal{N} (K_{fy} K_{yy}^{-1} (y(204) - m_y), K_{ff} - K_{fy} K_{yy}^{-1} K_{fy}^T)$$



An impractical but intuitively appealing way to think about Bayesian reasoning in the GP in the noiseless case:  
Of all the possible functions, just keep the ones that fit the data



A practical problem:  
choosing the best set of  
hyperparameters

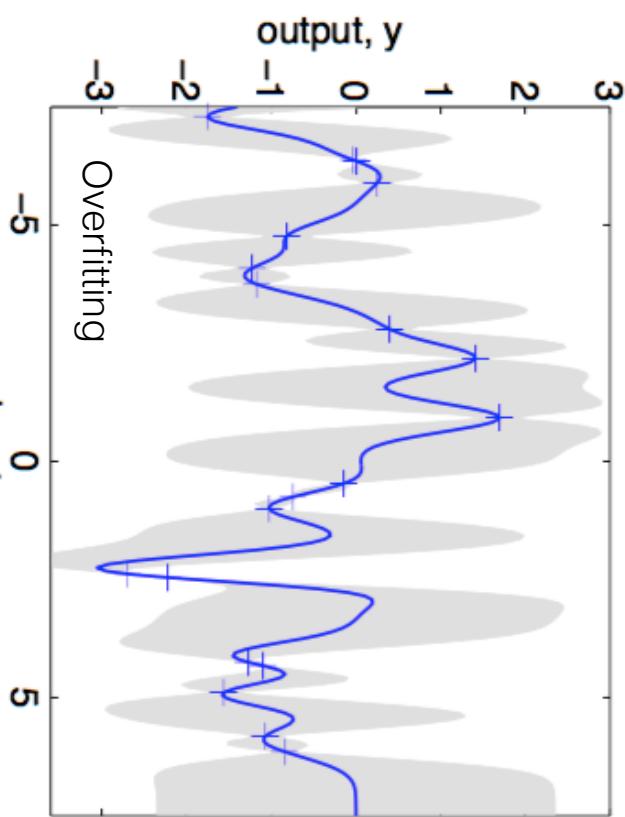
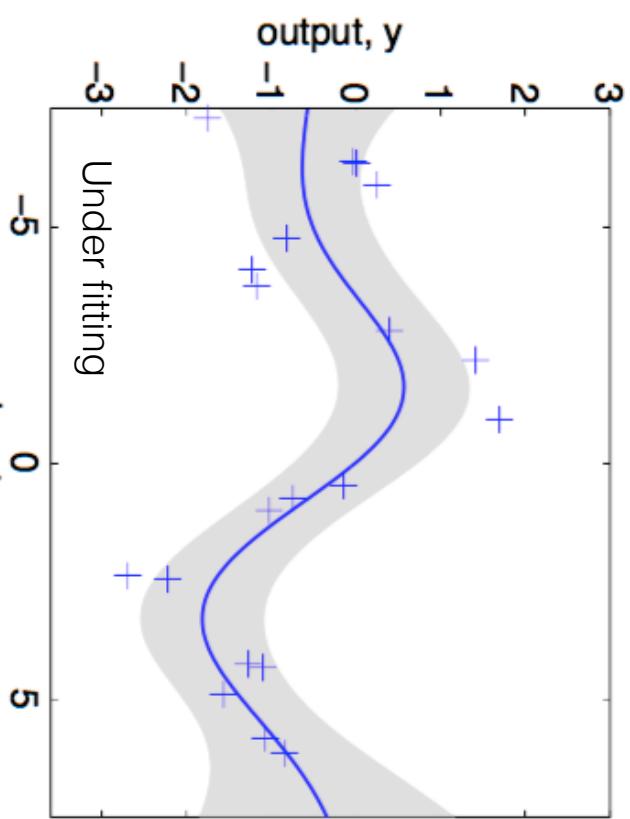
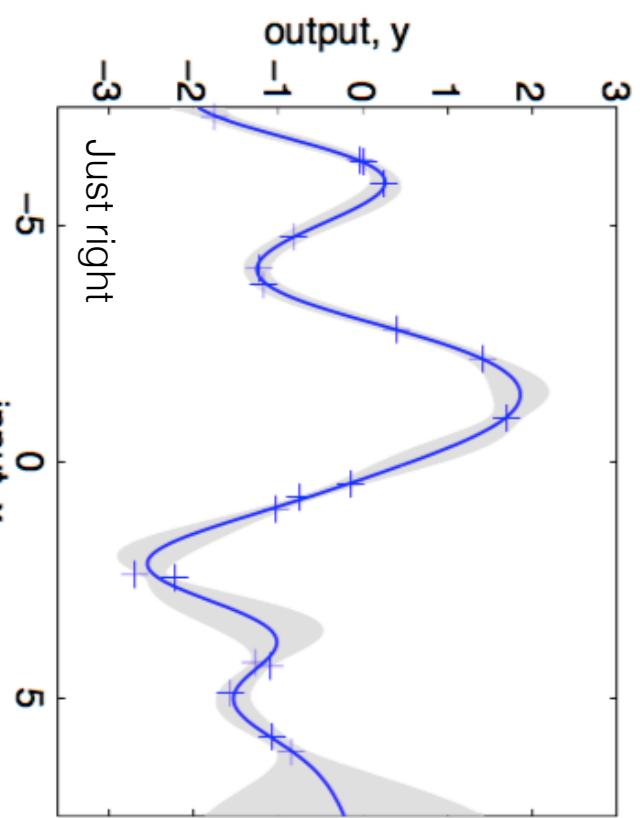
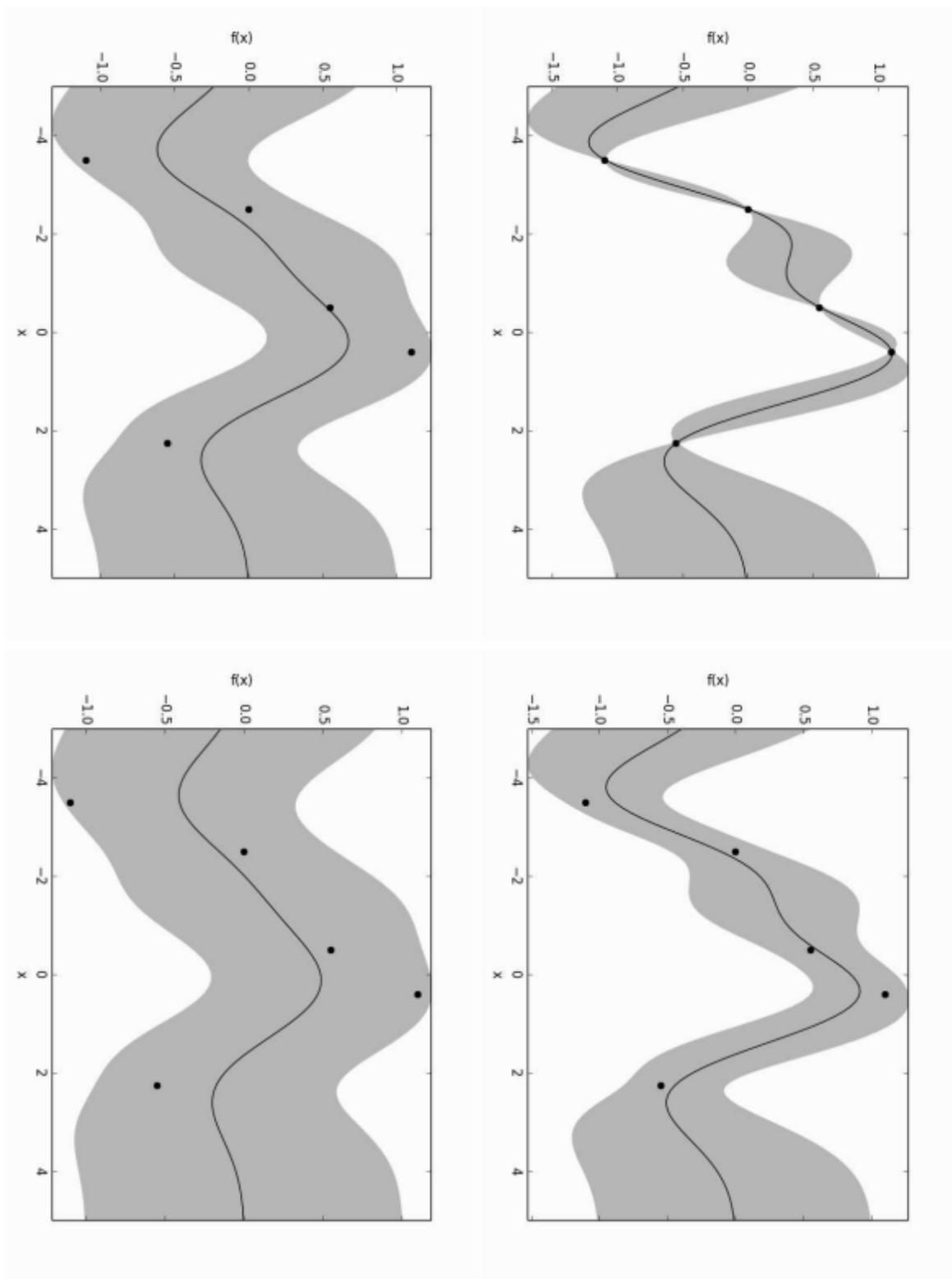


FIGURE 18.4 The effects of adding noise to the estimate of the covariance in the training data with the squared exponential kernel. Each plot shows the mean and 2 standard deviation error bars for a Gaussian process fitted to the five datapoints marked with dots. *Top left*:  $\sigma_n = 0.0$ , *top right*:  $\sigma_n = 0.2$ , *bottom left*:  $\sigma_n = 0.4$ , *bottom right*:  $\sigma_n = 0.6$ .



Now for learning hyper parameters ...

$$P(\mathbf{f}|\mathbf{X}) = (2\pi)^{\frac{n}{2}} |\mathbf{K}|^{\frac{1}{2}} \exp\left(-\frac{1}{2}\mathbf{f}^T \mathbf{K}^{-1} \mathbf{f}\right)$$

Simpler to work with log likelihood ...

$$\mathbf{f}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$$

$$\log p(\mathbf{f}|X) = -\frac{1}{2}\mathbf{f}^T K^{-1}\mathbf{f} - \frac{1}{2}\log|K| - \frac{n}{2}\log 2\pi.$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{0}, K + \sigma_n^2 I)$$

$$\log p(\mathbf{y}|X) = -\frac{1}{2}\mathbf{y}^T(K + \sigma_n^2 I)^{-1}\mathbf{y} - \frac{1}{2}\log|K + \sigma_n^2 I| - \frac{n}{2}\log 2\pi$$

Marsland's notation:  $\mathbf{t}$  instead of  $\mathbf{y}$

$$\log P(\mathbf{t}|\mathbf{x}, \boldsymbol{\theta}) = -\frac{1}{2}\mathbf{t}^T(\mathbf{K} + \sigma_n^2\mathbf{I})^{-1}\mathbf{t} - \frac{1}{2}\log|\mathbf{K} + \sigma_n^2\mathbf{I}| - \frac{N}{2}\log 2\pi$$

$$\mathbf{Q} = (\mathbf{K} + \sigma_n^2\mathbf{I})$$

Take derivatives wrt a parameter  $\theta$

$$\begin{aligned}\frac{\partial \mathbf{Q}^{-1}}{\partial \theta} &= -\mathbf{Q}^{-1}\frac{\partial \mathbf{Q}}{\partial \theta}\mathbf{Q}^{-1} \\ \frac{\partial \log|\mathbf{Q}|}{\partial \theta} &= \text{trace}\left(\mathbf{Q}^{-1}\frac{\partial \mathbf{Q}}{\partial \theta}\right)\end{aligned}$$

$$\begin{aligned}
k(\mathbf{x}, \mathbf{x}') &= \exp(\sigma_f) \exp\left(-\frac{1}{2} \exp(\sigma_l) |\mathbf{x} - \mathbf{x}'|^2\right) + \exp(\sigma_n) \mathbf{I} \\
&= k' + \exp(\sigma_n) \mathbf{I}
\end{aligned}$$

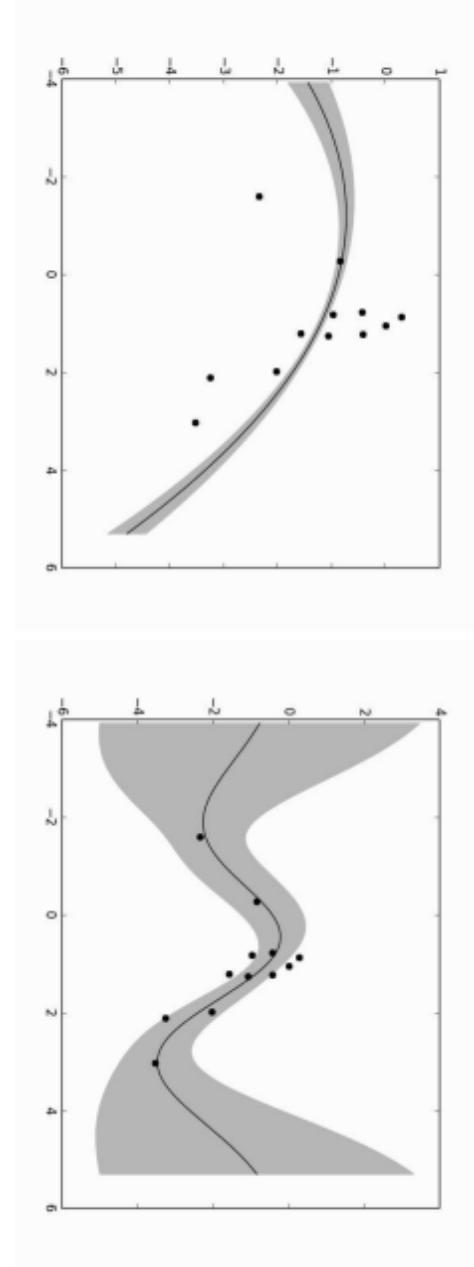
(Use  $\exp(\text{parameter})$  to make sure they are positive)

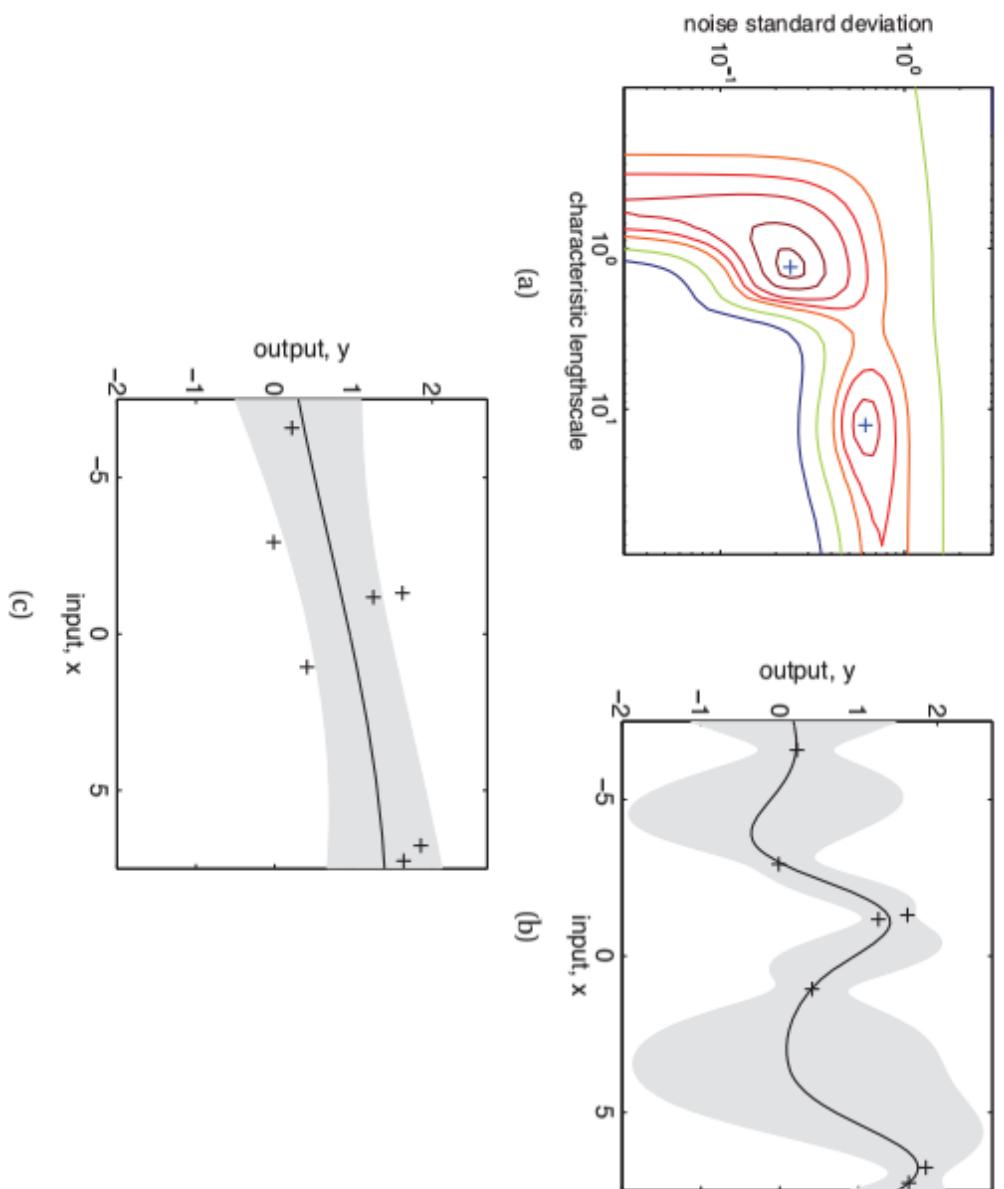
$$\frac{\partial k}{\partial \sigma_f} = k'$$

$$\frac{\partial k}{\partial \sigma_l} = k' \times \left( -\frac{1}{2} \exp(\sigma_l) |\mathbf{x} - \mathbf{x}'|^2 \right)$$

$$\frac{\partial k}{\partial \sigma_n} = \exp(\sigma_n) \mathbf{I}$$

Initial condition  
fitted parameters

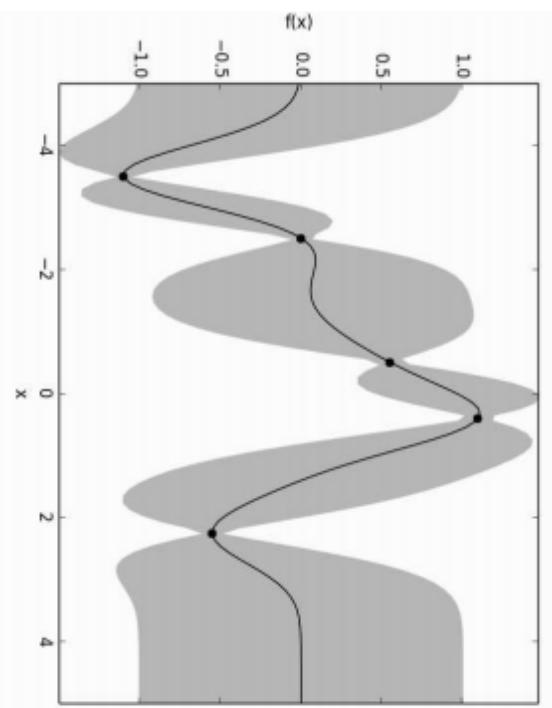




**Figure 15.5** Illustration of local minima in the marginal likelihood surface. (a) We plot the log marginal likelihood vs  $\sigma_y^2$  and  $\ell$ , for fixed  $\sigma_f^2 = 1$ , using the 7 data points shown in panels b and c. (b) The function corresponding to the lower left local minimum,  $(\ell, \sigma_n^2) \approx (1, 0.2)$ . This is quite “wiggly” and has low noise. (c) The function corresponding to the top right local minimum,  $(\ell, \sigma_n^2) \approx (10, 0.8)$ . This is quite smooth and has high noise. The data was generated using  $(\ell, \sigma_n^2) = (1, 0.1)$ . Source: Figure 5.5 of (Rasmussen and Williams 2006). Figure generated by gprDemoMarglik, written by Carl Rasmussen.

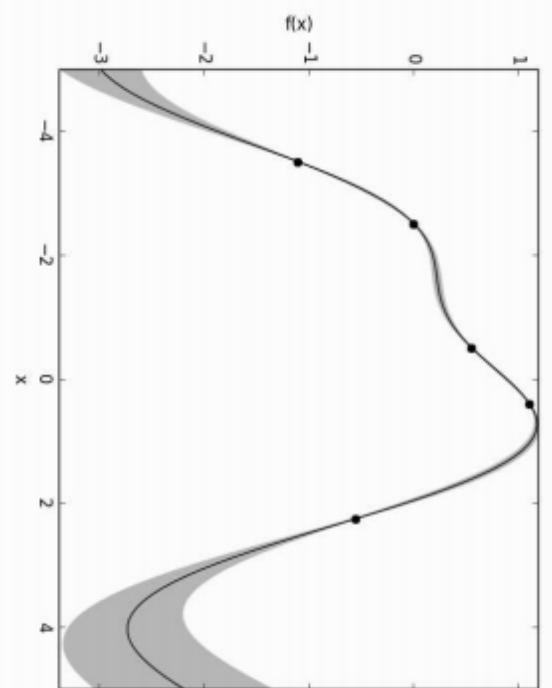
$X$ s are relatively uncorrelated

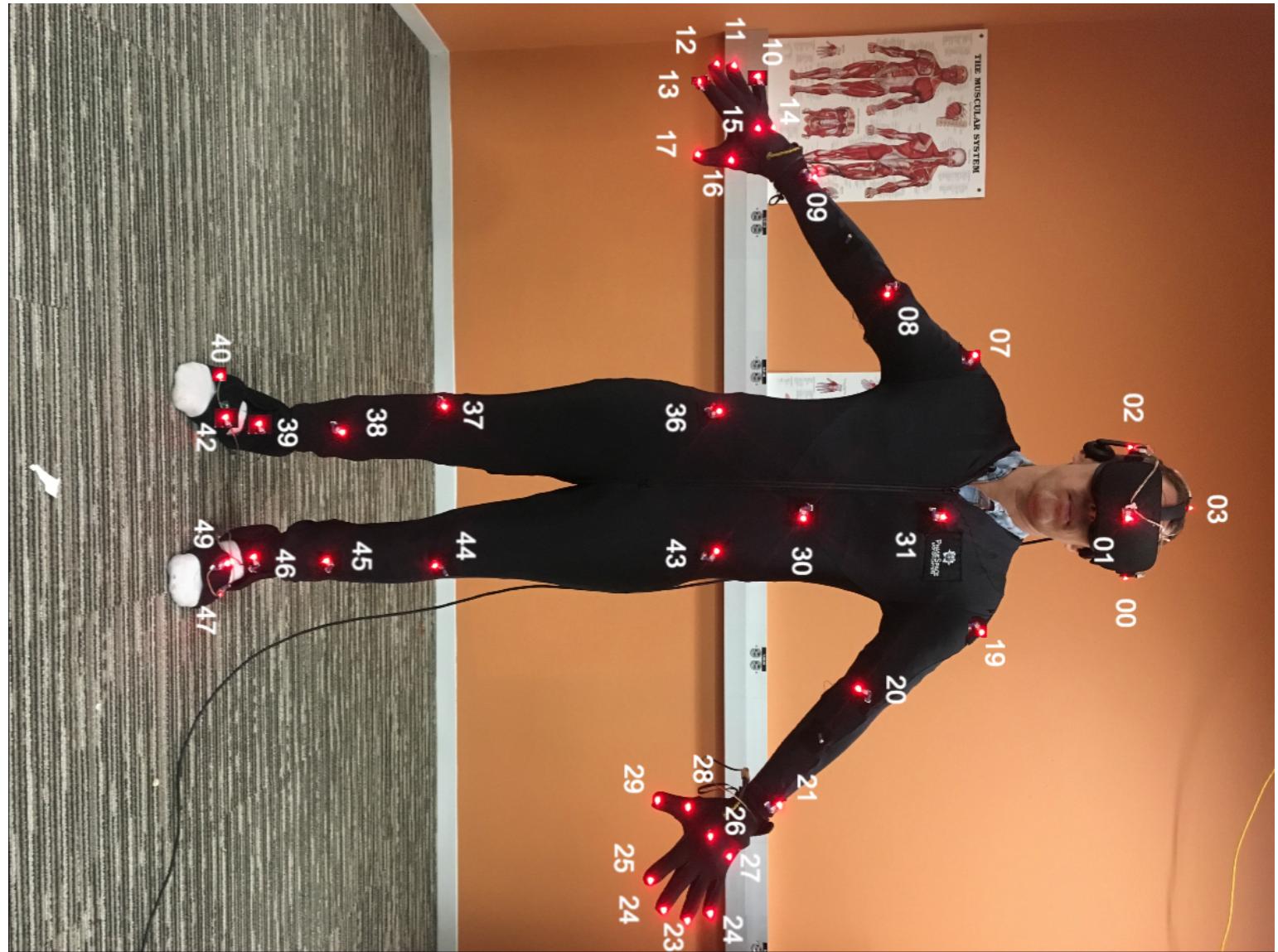
$$|=0.5$$

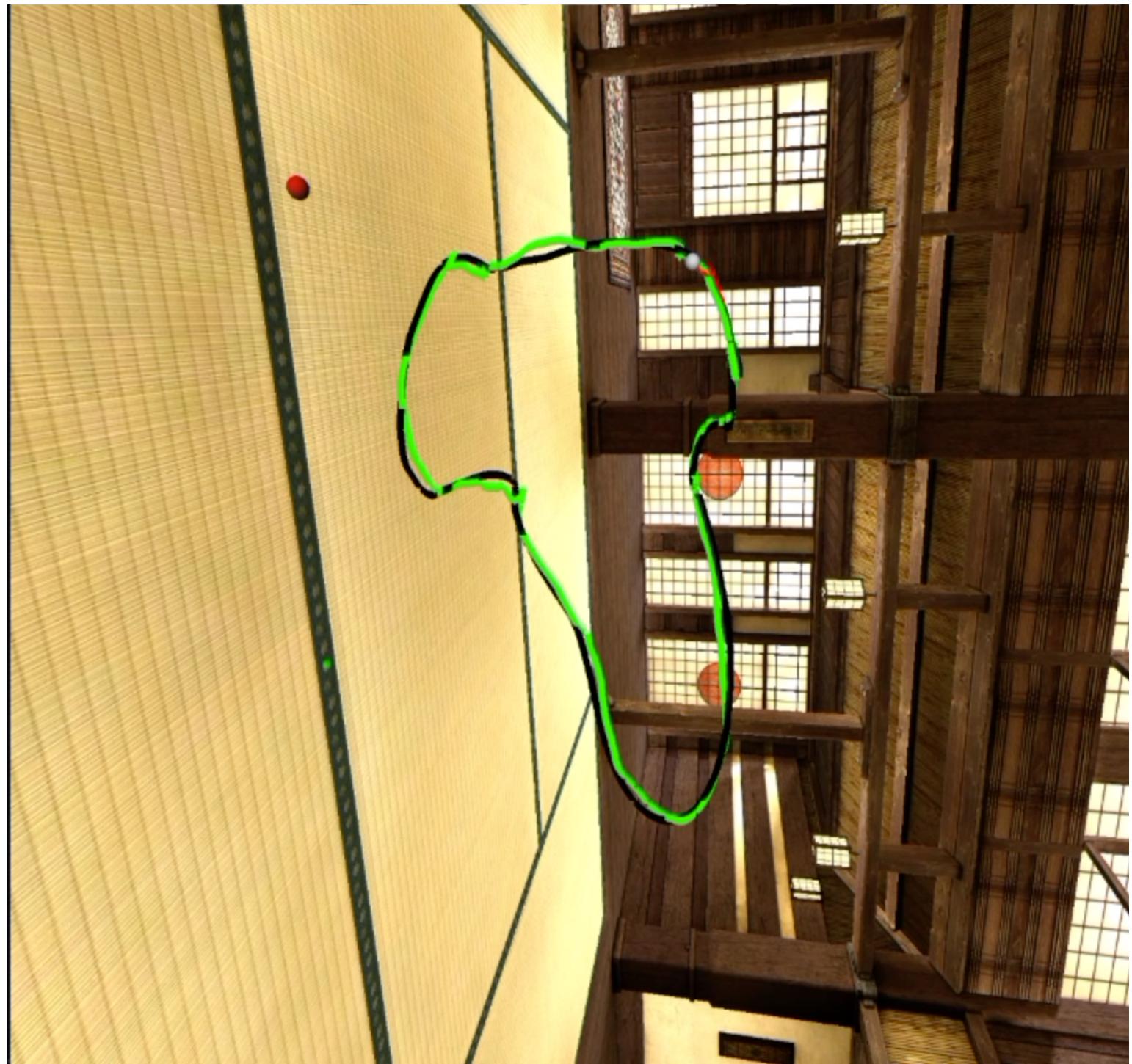


$X$ s are more uncorrelated

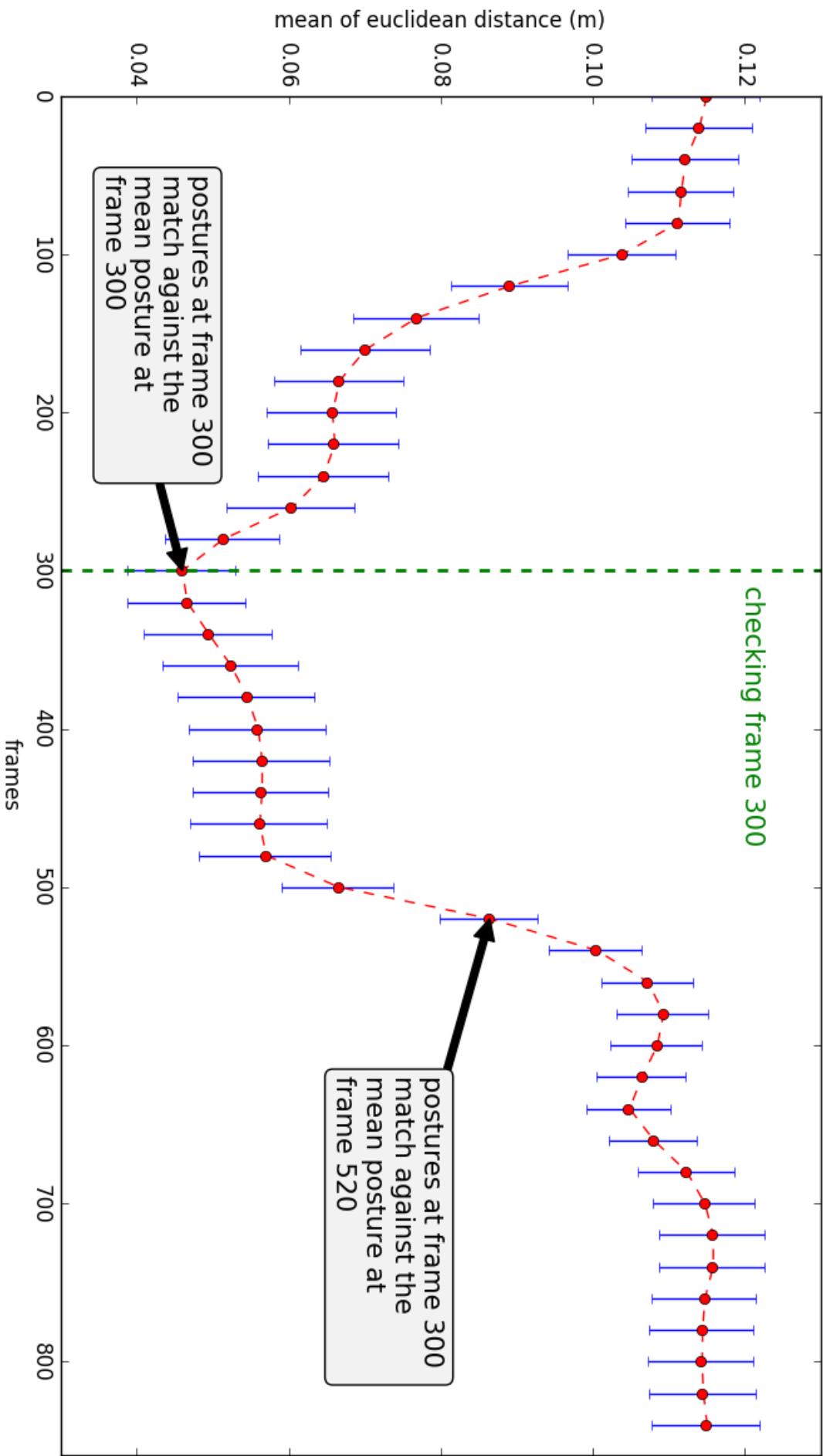
$$|=2.0$$







The postures at frame 300 do not look like postures at any other frame



# Dynamics: Notation

Symbol	Meaning
$x$	position or state of one or more rigid bodies
$\dot{x}$	velocity (usu. linear and angular)
$\ddot{x}$	acceleration (usu. linear and angular)
$R$	rotation matrix representing orientation of a body
$\omega$	angular velocity
$q$	quaternion representation of an orientation or rotation
$m$	mass of a single rigid body
$M$	mass matrix
$I$	identity matrix
$\mathcal{I}$	moment of inertia tensor
$n_b, n_c$	number of bodies, number of constraints
$\alpha, \beta$	stabilizing parameters added to the equations of motion
$\phi()$	error or energy function for a single constraint
$J$	matrix of partial derivatives of constraint error functions
$h$	timestep
$f$	forces (and torques)
$\tau$	torques
$\lambda$	constraint forces

Table 1: Meanings of specific symbols used to discuss dynamic simulation

## Newton's Laws

$$\begin{bmatrix} f_{1t} \\ \tau_{1t} \\ f_{2t} \\ \tau_{2t} \end{bmatrix} = \begin{bmatrix} m_1 I & 0 & 0 & 0 \\ 0 & \mathcal{I}_{1t} & 0 & 0 \\ 0 & 0 & m_2 I & 0 \\ 0 & 0 & 0 & \mathcal{I}_{2t} \end{bmatrix} \begin{bmatrix} \ddot{x}_{1t} \\ \dot{\omega}_{1t} \\ \ddot{x}_{2t} \\ \dot{\omega}_{2t} \end{bmatrix} \Rightarrow f_t = M_t \ddot{x}_t$$

Forces have gyroscopic, gravity and control components

$$f = f_c + f_g + f_u.$$

Semi-implicit Euler uses future velocity for computing position

$$\dot{x}_{t+h} = \dot{x}_t + h M^{-1} f_t \quad (1)$$

$$x_{t+h} = x_t + h \dot{x}_{t+h} \quad (2)$$

# Handling Joint Constraints

$$J_t = \nabla \phi(x_t) = \begin{bmatrix} \frac{\partial \phi_1}{\partial x_{1t}} & \dots & \frac{\partial \phi_1}{\partial x_{(6n_b)t}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \phi_k}{\partial x_{1t}} & \dots & \frac{\partial \phi_k}{\partial x_{(6n_b)t}} \end{bmatrix}$$

Reduce joint constraints gradually during the integration

$$(I - \alpha)\phi(x_t) = \phi(x_{t+h}) \approx \phi(x_t) + J_t(x_{t+h} - x_t)$$

$$\dot{\mathbf{x}}_{t+h} = \dot{\mathbf{x}}_t + h \mathbf{M}^{-1} \mathbf{f}_t \quad (1)$$

$$\mathbf{x}_{t+h} = \mathbf{x}_t + h \dot{\mathbf{x}}_{t+h}$$

$$(I - \alpha) \phi(\mathbf{x}_t) = \phi(\mathbf{x}_{t+h}) \approx \phi(\mathbf{x}_t) + \mathbf{J}_t (\mathbf{x}_{t+h} - \mathbf{x}_t) \quad (3)$$

$$\mathbf{x}_{t+h} = \mathbf{x}_t + h \dot{\mathbf{x}}_t + h^2 \mathbf{M}_t^{-1} (\mathbf{f}_{ct} + \mathbf{f}_{gt} + \mathbf{f}_{ut}) \quad (4)$$

Use 1~4 to eliminate all future quantities t+h

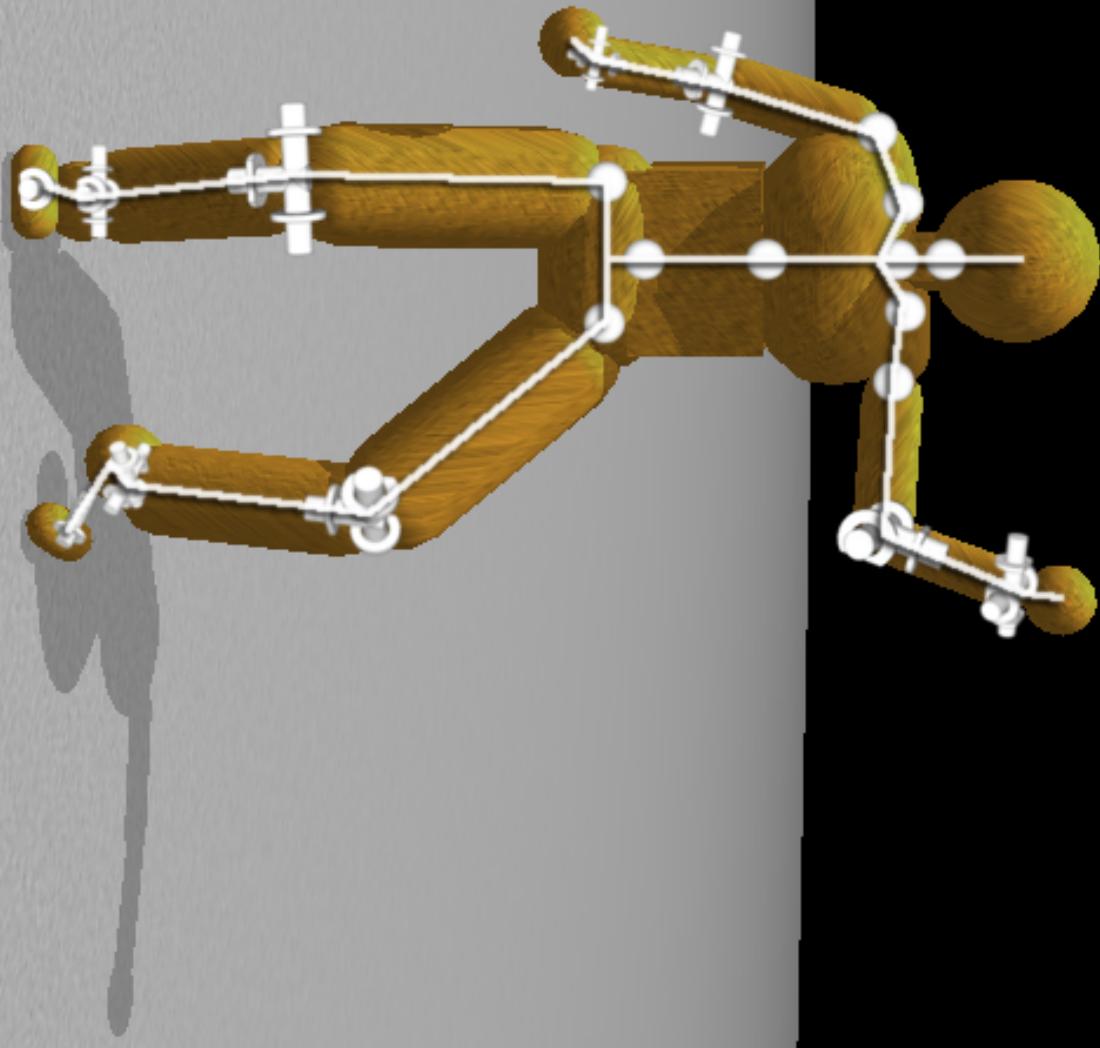
$$(I - \alpha) \phi(\mathbf{x}_t) = \phi(\mathbf{x}_t) + \mathbf{J}_t (\mathbf{x}_t + h \dot{\mathbf{x}}_t + h^2 \mathbf{M}_t^{-1} (\mathbf{f}_{ct} + \mathbf{f}_{gt} + \mathbf{f}_{ut}) - \mathbf{x}_t)$$

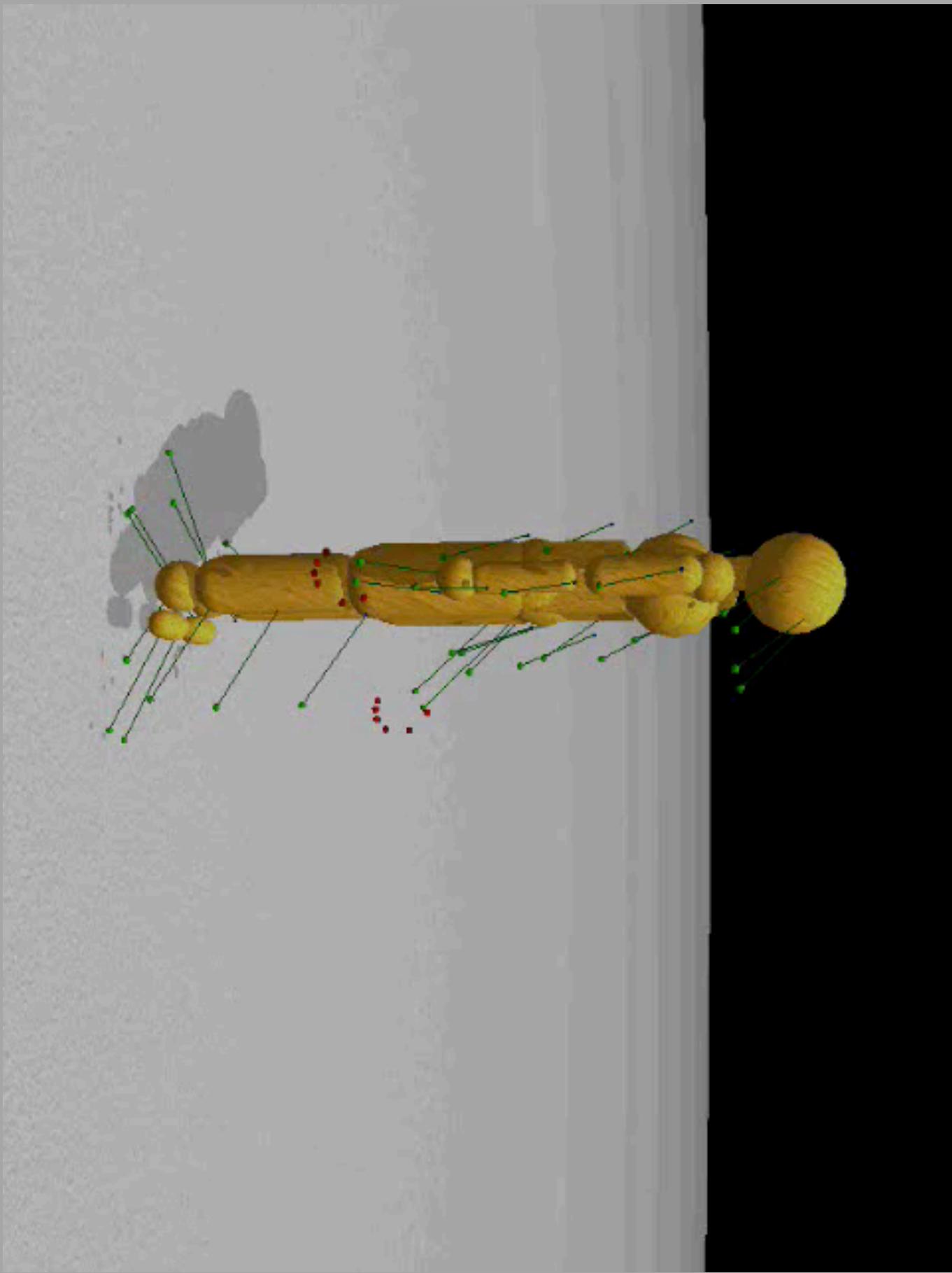
Solve for  $\mathbf{f}_{ct}$

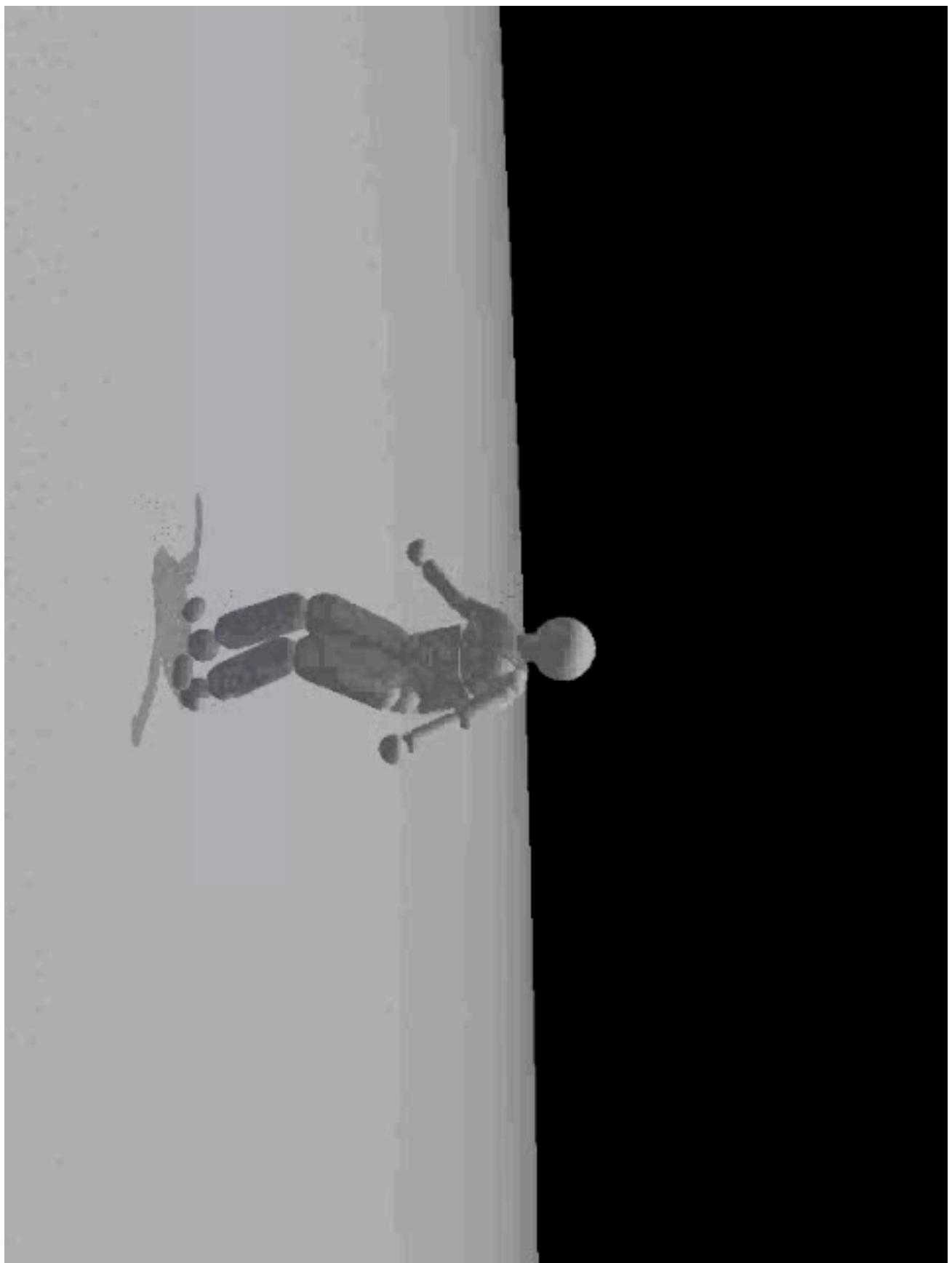
$$\mathbf{J}_t \mathbf{M}_t^{-1} \mathbf{f}_{ct} = -\frac{1}{h^2} \alpha \phi(\mathbf{x}_t) - \frac{1}{h} \mathbf{J}_t \dot{\mathbf{x}}_t - \mathbf{J}_t \mathbf{M}_t^{-1} (\mathbf{f}_{gt} + \mathbf{f}_{ut})$$

## Dynamic Model

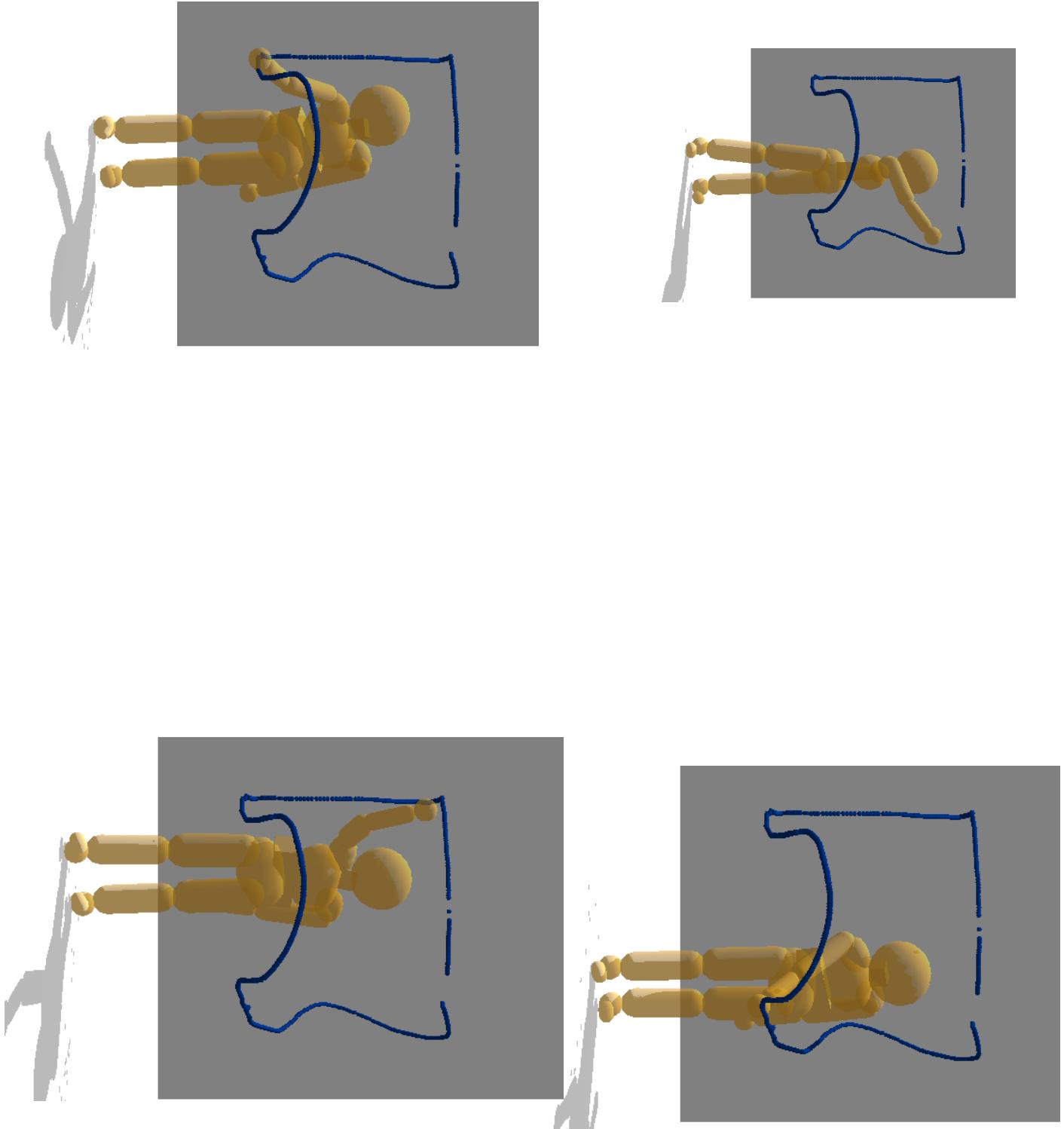
Including  
50 DOFs,  
Lengths,  
Masses,  
Inertial Tensors











## One subject traces the square curve

