

Causal estimation of 3D structure and motion

Stefano Soatto
UCLA

VISION as a **SENSOR**

machine to **INTERACT** with the environment

NEED to estimate relative 3D **MOTION**

3D **SHAPE** (TASK)

REAL-TIME

CAUSAL processing

representation of **SHAPE**

(only supportive of representation of motion)

POINT-FEATURES

TRADEOFFS

SFM

CORRESPONDENCE

LARGE
BASELINE

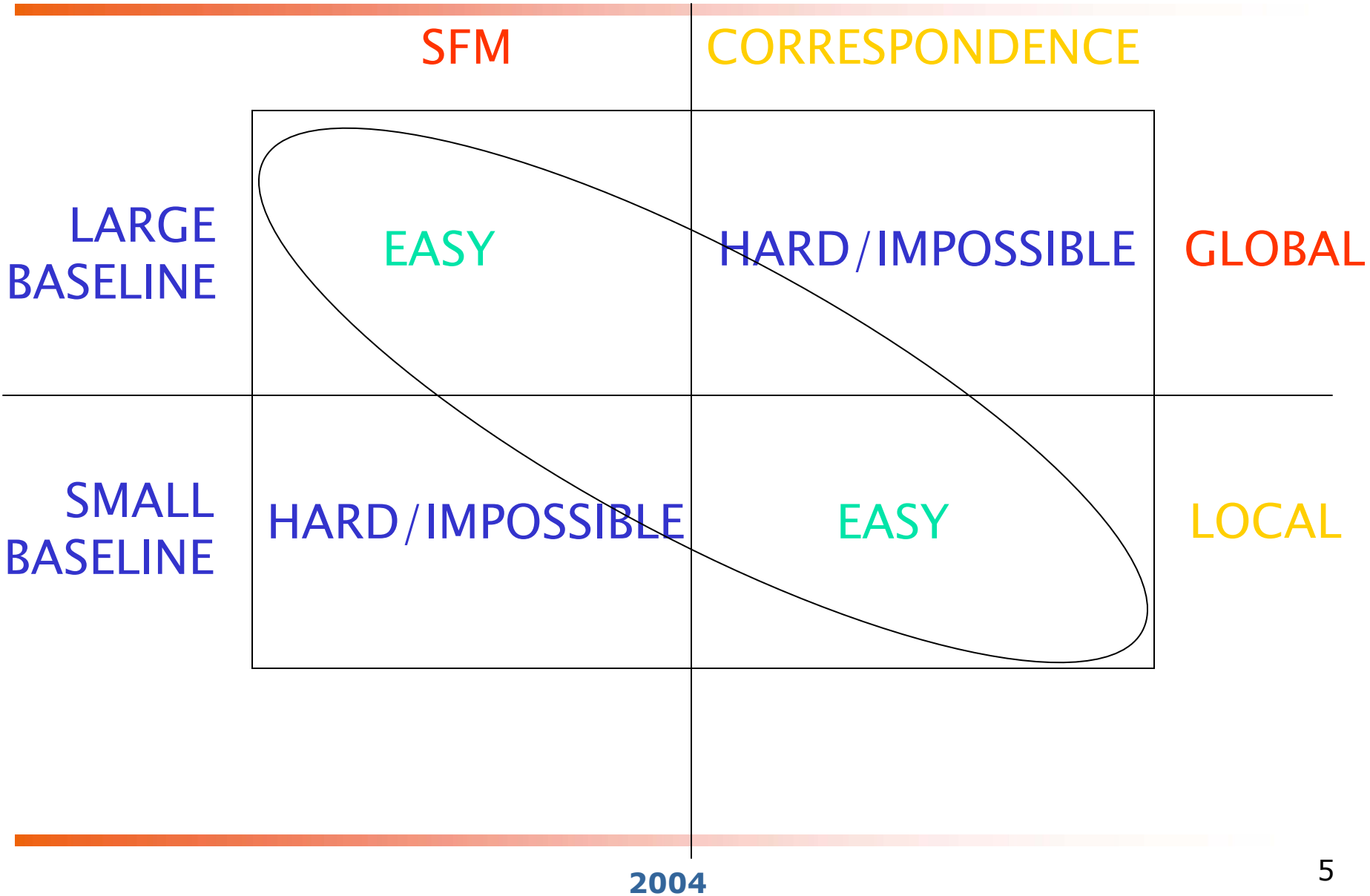
EASY

HARD/IMPOSSIBLE

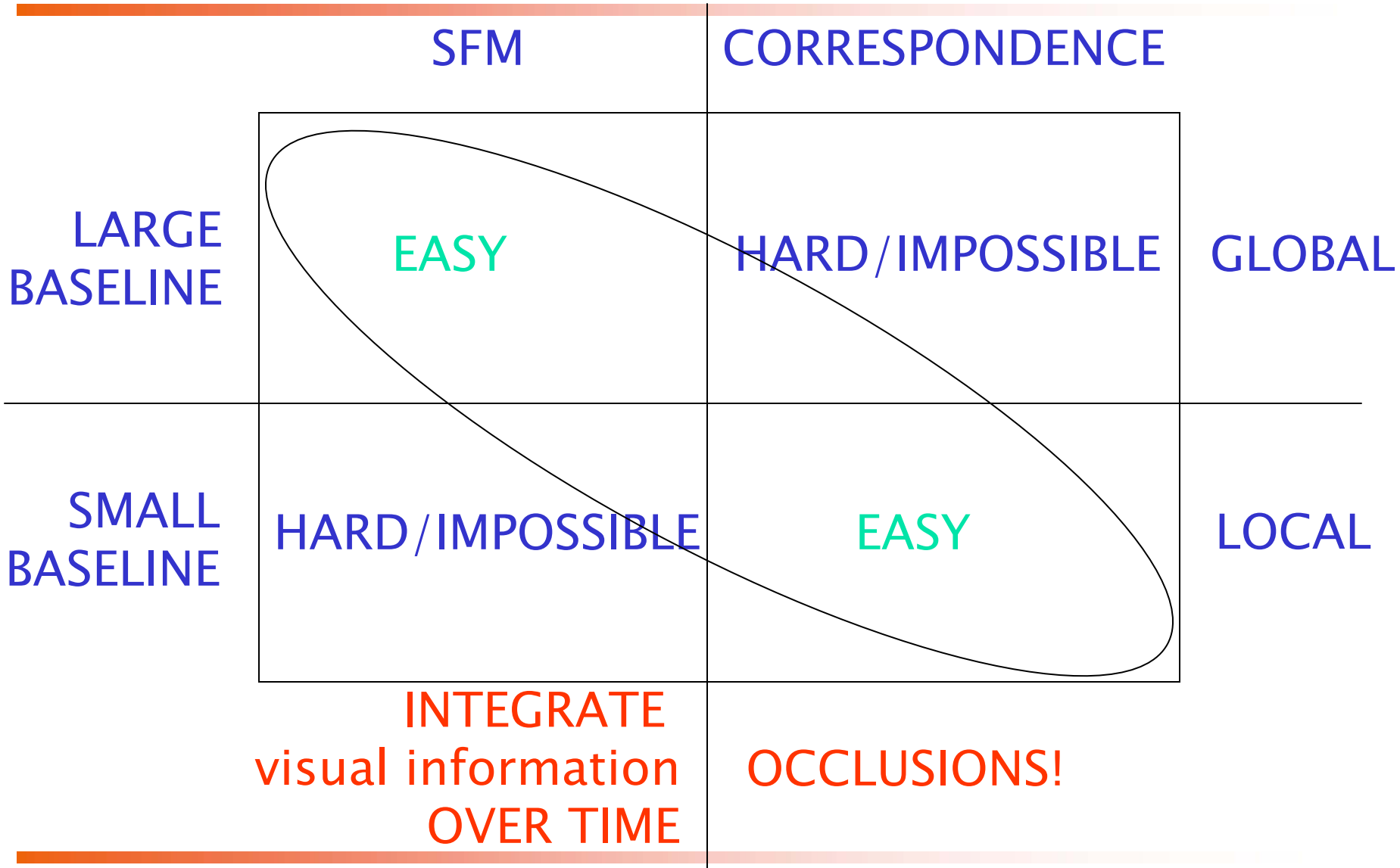
TRADEOFFS

	SFM	CORRESPONDENCE
LARGE BASELINE	EASY	HARD/IMPOSSIBLE
SMALL BASELINE	HARD/IMPOSSIBLE	EASY

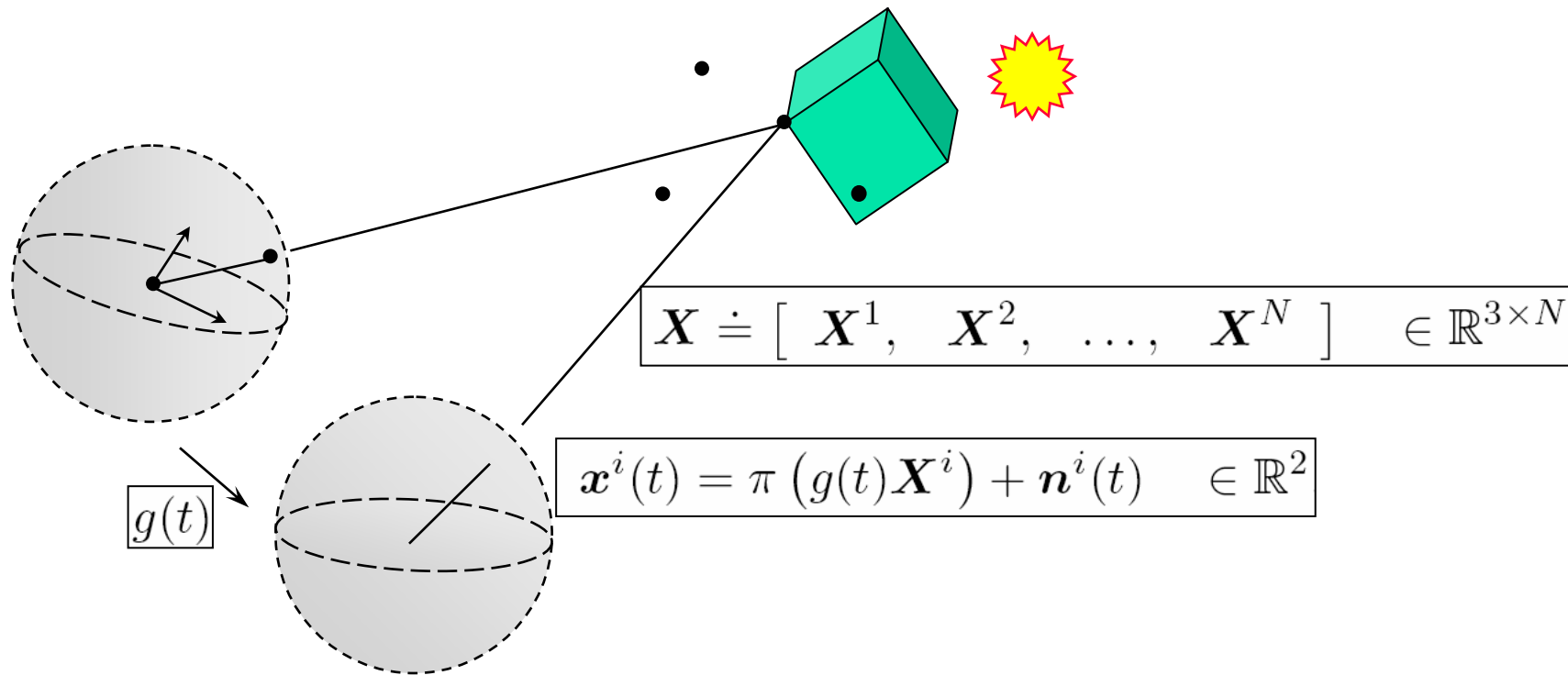
TRADEOFFS



WHAT DOES IT TAKE ?



Setup and notation



Temporal evolution:

$$\mathbf{X}(t+1) = \mathbf{X}(t)$$

$$g(t+1) = \exp(\widehat{\xi}(t))g(t)$$

Structure and motion as a filtering problem

$$\left\{ \begin{array}{ll} \mathbf{X}(t+1) = \mathbf{X}(t), & \mathbf{X}(0) = \mathbf{X}_0 \in \mathbb{R}^{3 \times N}, \\ g(t+1) = \exp(\widehat{\xi}(t))g(t), & g(0) = g_0 \in SE(3), \\ \xi(t+1) = \xi(t) + \alpha(t), & \xi(0) = \xi_0 \in \mathbb{R}^6, \\ \mathbf{x}^i(t) = \pi(g(t)\mathbf{X}^i(t)) + \mathbf{n}^i(t), & \mathbf{n}^i(t) \sim \mathcal{N}(0, \Sigma_n). \end{array} \right.$$

- Given measurements of the “output” (feature point positions)
- Given modeling assumptions about the “input” (acceleration = noise)
- Estimate the “state” (3D structure, pose, velocity)

Difficulties ...

$$\left\{ \begin{array}{ll} \mathbf{X}(t+1) = \mathbf{X}(t), & \mathbf{X}(0) = \mathbf{X}_0 \in \mathbb{R}^{3 \times N}, \\ g(t+1) = \exp(\widehat{\xi}(t))g(t), & g(0) = g_0 \in SE(3), \\ \xi(t+1) = \xi(t) + \alpha(t), & \xi(0) = \xi_0 \in \mathbb{R}^6, \\ \mathbf{x}^i(t) = \pi(g(t)\mathbf{X}^i(t)) + \mathbf{n}^i(t), & \mathbf{n}^i(t) \sim \mathcal{N}(0, \Sigma_n). \end{array} \right.$$

- Model is non-linear (output map = projection)
- State-space is non-linear! (SE(3))
- Noise: need to specify what we mean by “estimate”
- Even without noise: *model is not observable!*

Observability

- Equivalent class of state-space trajectories generate the same measurements

$$\{\mathbf{X}_0, g_0, \xi_0\} \quad g_0 = (R_0, T_0) \quad \xi_0 = (\omega_0, v_0)$$

$$\{\beta \tilde{R} \mathbf{X}_0 + \beta \tilde{T}, \tilde{g}_0, \tilde{\xi}_0\} \quad \tilde{g}_0 = (R_0 \tilde{R}^T, \beta T_0 - \beta R_0 \tilde{R}^T \tilde{T}) \quad \tilde{\xi}_0 = (\omega_0, \beta v_0)$$

- Fix, e.g., the direction of 3 points, and the depth of one point (Gauge transformation)

Local coordinatization of the state space

$$\left\{ \begin{array}{ll} \mathbf{x}_0^i(t+1) = \mathbf{x}_0^i(t), & i = 1, 2, \dots, N, & \mathbf{x}_0^i(0) = \mathbf{x}_0^i, \\ \lambda^i(t+1) = \lambda^i(t), & i = 1, 2, \dots, N, & \lambda^i(0) = \lambda_0^i, \\ T(t+1) = \exp(\widehat{\omega}(t))T(t) + v(t), & & T(0) = T_0, \\ \Omega(t+1) = \log_{SO(3)}(\exp(\widehat{\omega}(t))\exp(\widehat{\Omega}(t))), & & \Omega(0) = \Omega_0, \\ v(t+1) = v(t) + \alpha_v(t), & & v(0) = v_0, \\ \omega(t+1) = \omega(t) + \alpha_\omega(t), & & \omega(0) = \omega_0, \\ \mathbf{x}^i(t) = \pi\left(\exp(\widehat{\Omega}(t))\mathbf{x}_0^i(t)\lambda^i(t) + T(t)\right) + \mathbf{n}^i(t), & i = 1, 2, \dots, N. & \end{array} \right.$$

$\log_{SO(3)}(R)$ stands for Ω such that $R = e^{\widehat{\Omega}}$

Minimal realization

$$\left\{ \begin{array}{ll} \mathbf{x}_0^i(t+1) = \mathbf{x}_0^i(t), & i = 4, 5, \dots, N, & \mathbf{x}_0^i(0) = \mathbf{x}_0^i, \\ \lambda^i(t+1) = \lambda^i(t), & i = 2, 3, \dots, N, & \lambda^i(0) = \lambda_0^i, \\ T(t+1) = \exp(\widehat{\omega}(t))T(t) + v(t), & & T(0) = T_0, \\ \Omega(t+1) = \log_{SO(3)}(\exp(\widehat{\omega}(t))\exp(\widehat{\Omega}(t))), & & \Omega(0) = \Omega_0, \\ v(t+1) = v(t) + \alpha_v(t), & & v(0) = v_0, \\ \omega(t+1) = \omega(t) + \alpha_\omega(t), & & \omega(0) = \omega_0, \\ \mathbf{x}^i(t) = \pi\left(\exp(\widehat{\Omega}(t))\mathbf{x}_0^i(t)\lambda^i(t) + T(t)\right) + \mathbf{n}^i(t), & i = 1, 2, \dots, N. \end{array} \right.$$

$$x(t) \doteq [\mathbf{x}_0^4(t)^T, \dots, \mathbf{x}_0^N(t)^T, \lambda^2(t), \dots, \lambda^N(t), T^T(t), \Omega^T(t), v^T(t), \omega^T(t)]^T,$$

$$y(t) \doteq [\mathbf{x}^1(t)^T, \dots, \mathbf{x}^N(t)^T]^T.$$

- Now it looks very much like: $\begin{cases} x(t+1) = f(x(t)), & x(t_0) = x_0 \\ y(t) = h(x(t)), \end{cases}$
- And we are looking for (a point statistic of): $p(x(t)|y^t)$

$$\hat{x}(t|t) = E_p[x(t)] = \int x(t) dP(x(t)|y^t)$$

EKF vs. particle filter?

- For single rigid body/static scene, expect unimodal posterior
- No need to estimate entire density; point estimate suffices
- Robust (M-) version of EKF works well in practice ...
- ... and in real time for a few hundred feature points

In practice ...

- Adding/removing features (subfilters)
- Multiple motions/outliers (M-filter, innovation tests)
- Tracking drift (reset with wide-baseline matching)
- Switching the reference features (hard! Causes unavoidable global drift)
- Global registration (maintain DB of lost features)

QuickTime™ and a
MS-MPEG4v2 Video decompressor
are needed to see this picture.