

# Estimating 3D Motion and Forces of Person-Object Interactions From Monocular Video

*Authors: Zongmian Li, Jiri Sedlar, Justin Carpentier, Ivan Laptev, Nicolas Mansard, Josef Sivic  
(CVPR-2019)*

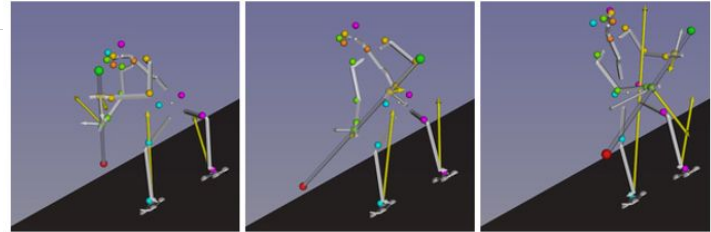
Sahana Eshwaran, Alicia Alarie, Simba Nyatsanga

# Problem Statement

Input: A monocular RGB video



Output 3D Motion and Contact Forces



Given an input video, reconstruct 3D motion and forces of a person interacting with an object from a single RGB video

# Motivation

- Large-scale learning of human-object interactions from internet videos, speed up the pace of learning
- Advanced visual intelligence capabilities to recognize and interpret complex person-to-object interactions

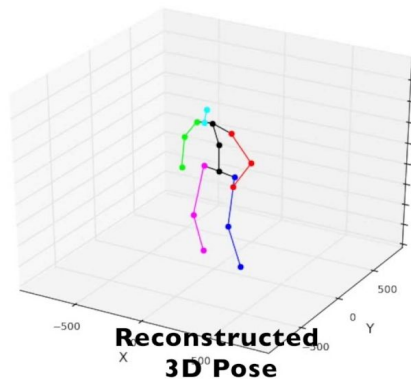
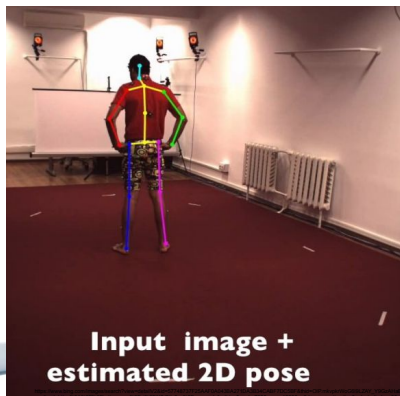


"Atlas" by Boston Dynamics

UC DAVIS

# Related Work

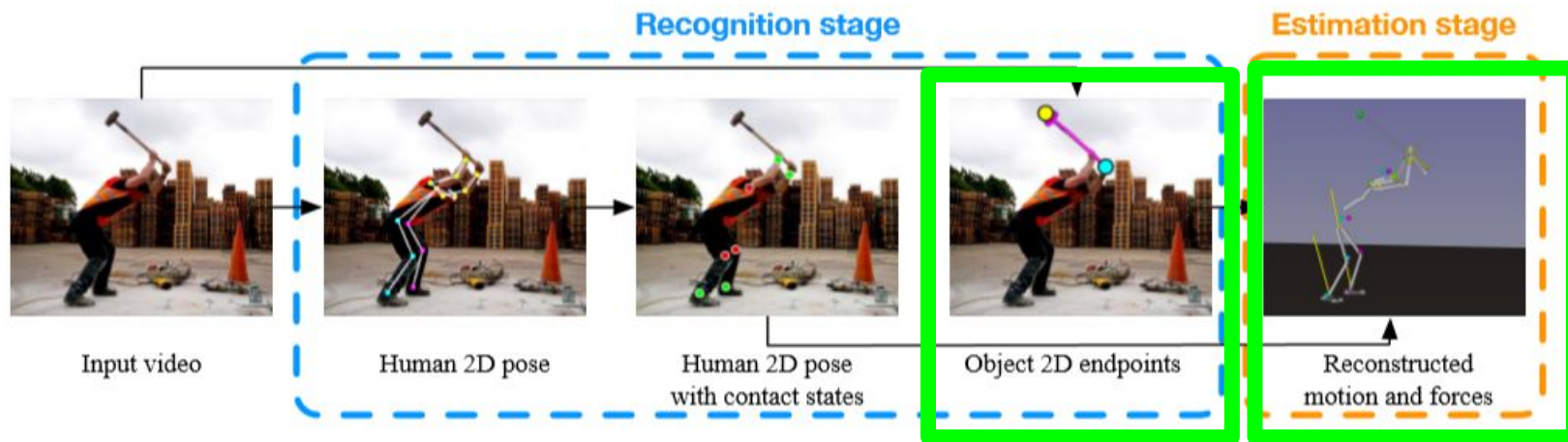
- Single-View 3D pose estimation
- Human- Object Interactions
- Object 3D pose estimation
- Learning from instructional videos



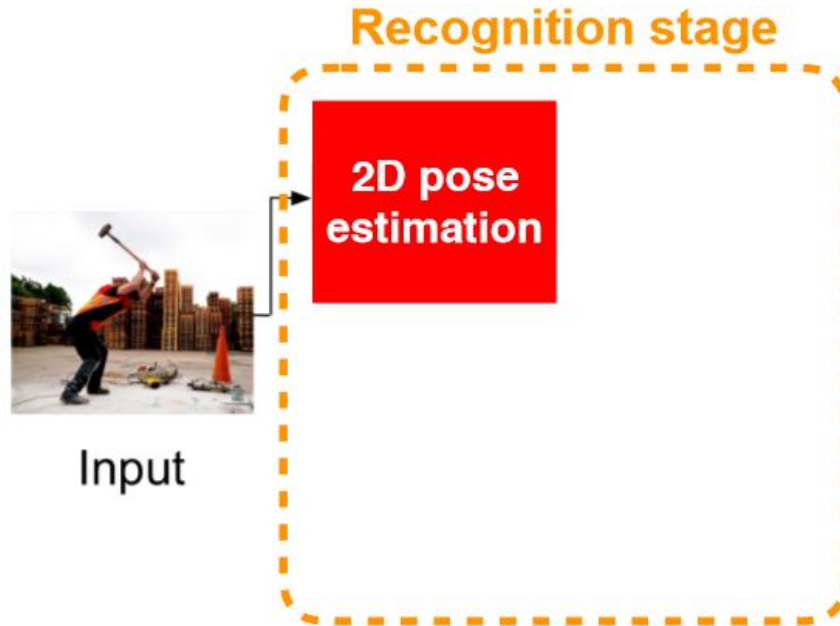
# Work borrowed from Person-object interactions in robotics.

- They use the concept of rigid body models introduced in robotics and estimate the 3D person-object interactions from monocular videos using optimal control problem under contact constraints.
- Then identify the contact irregularities by identifying contact states from visual input & localize the contact points in 3D via trajectory estimator.

# Key Contributions



# Method: A two-stage approach



[Cao et al., CVPR 2017]



# Method: A two-stage approach

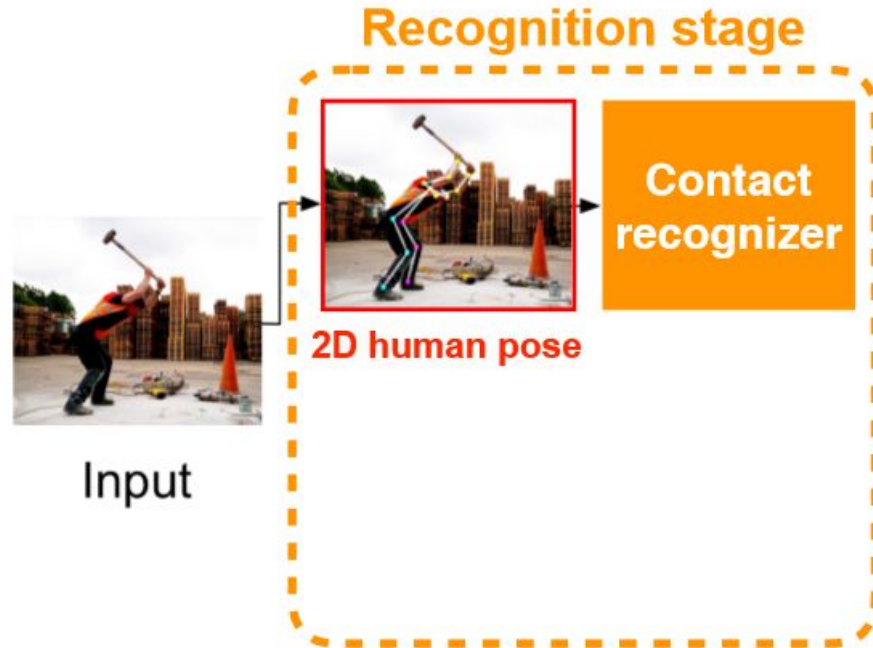


**2D pose  
estimation**



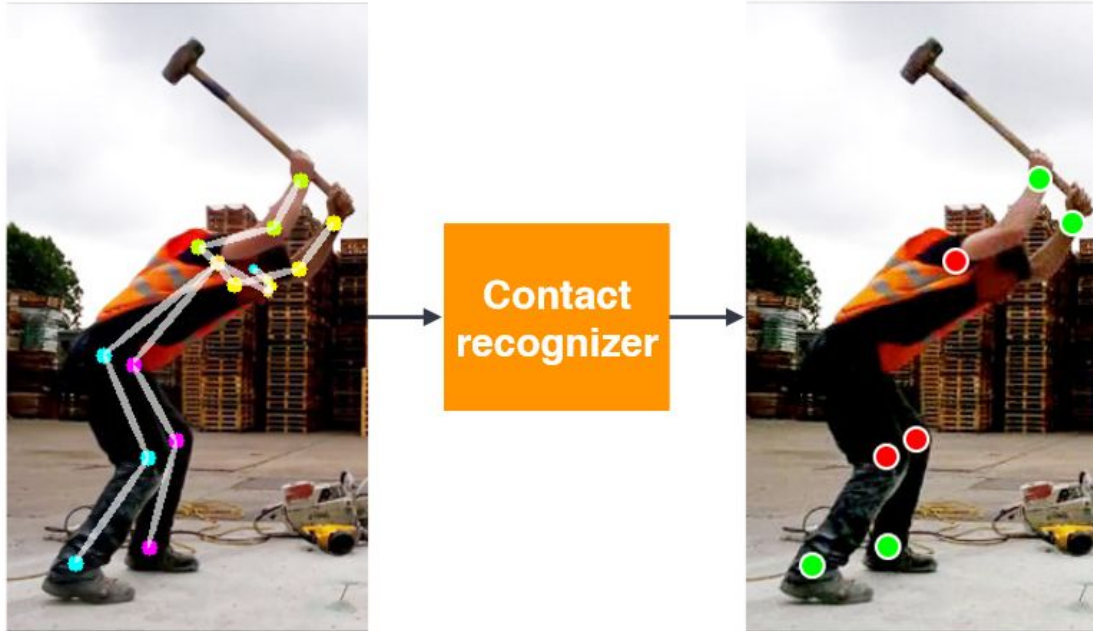


# Method: A two-stage approach

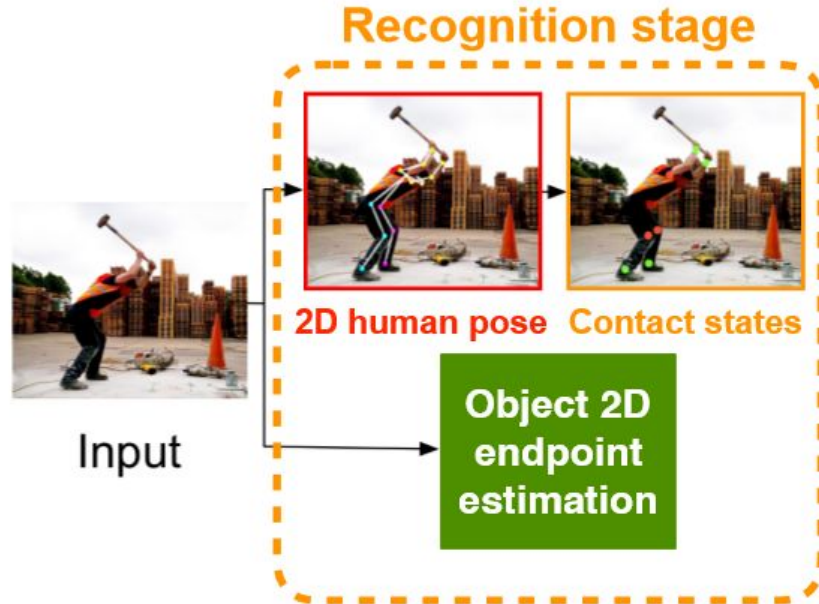


[He et al., CVPR 2016]

# Method: A two-stage approach

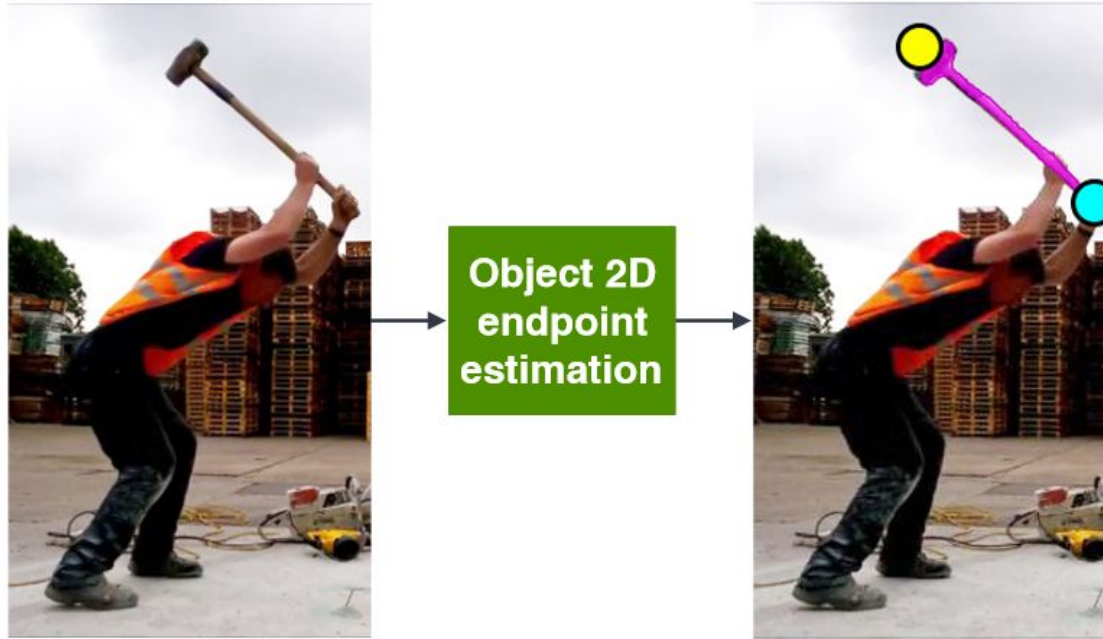


# Method: a two-stage approach

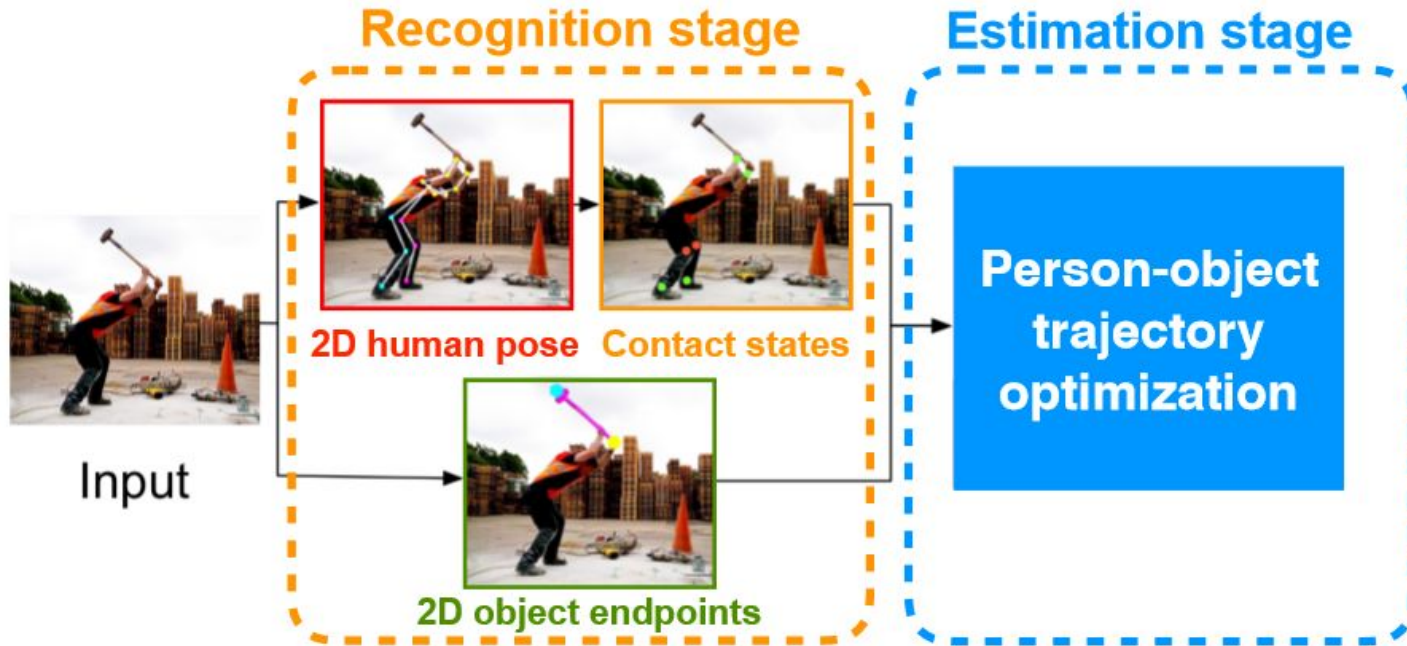


[He et al., ICCV 2017]

# Method: a two-stage approach



# Method: a two-stage approach



[Tassa et al., IROS 2012]  
[Carpentier et al., TRO, 2018]

# Estimation Stage

- Assume a video clip of duration  $T$  depicting a human manipulating an object.
- Encode the 3D poses of the human and object using configuration vectors (translation and rotational displacement).

$$\mathbf{x} = (q^h, q^o, \dot{q}^h, \dot{q}^o)$$

- $K$  is constant set of possible contact points between human, object and/or ground plane.

# Estimation Stage

- Model the contact force exerted at a contact point (e.g. grip, Newton's 3rd law)

$$f_k \text{ where } k = 1, \dots, K$$

- Model the joint torque vector describing actuation by human muscles

$$\tau_m^h$$

- Control variable combining joint torque vector and contact forces at the K contact joints

$$u = (\tau_m^h, f_k, k = 1, \dots, K)$$



# Estimation Stage

- For sliding contacts, define a constant  $c$  consisting of relative positions of all contact points w.r.t the object/ground
- **Objective 1** - estimate smooth and consistent human-object and contact trajectories:

$\underline{x}$                        $\underline{c}$

- **Objective 2** - recover the control value  $\underline{u}$  that causes the observed motion
- **Intuition** - human and object 3D poses = respective projections in the image

# Estimation Stage

- Formulate person-object interaction estimation as a optimal control problem with contact and dynamic constraints:

$$\underset{\underline{x}, \underline{u}, \underline{c}}{\text{minimise}} \sum_{e \in \{h, o\}} \int_0^T l^e(x, u, c) dt$$

- Constrained by

$$\kappa(x, c) = 0$$

$$\dot{x} = f(x, c, u)$$

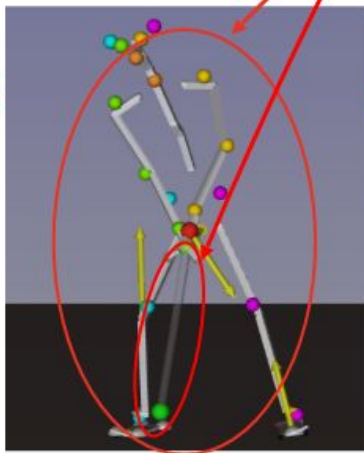
$$u \in \mathcal{U}$$

# Estimation Stage

$$\underset{\mathbf{x}, \underline{u}, \underline{c}}{\text{minimize}} \quad \sum_{e \in \{h, o\}} \int_0^T l^e(\mathbf{x}, u, c) dt$$

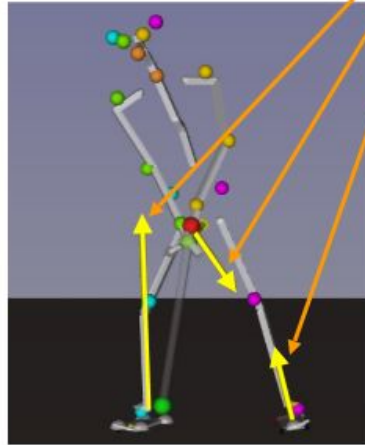
Person 3D poses

Object 3D poses



# Estimation Stage

$$\underset{\underline{x}, \underline{u}, \underline{c}}{\text{minimize}} \quad \sum_{e \in \{h, o\}} \int_0^T l^e(x, \underline{u}, c) dt$$

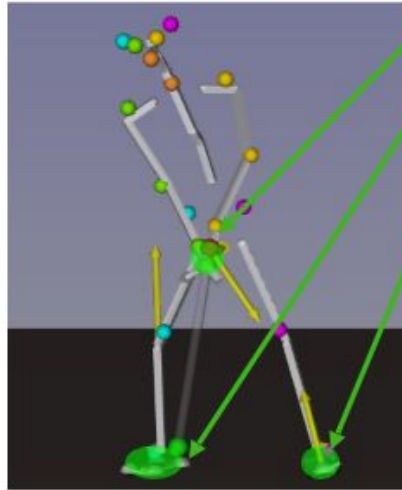


**Person-object  
person-ground  
contact forces**

# Estimation Stage

minimize  
 $\underline{x}, \underline{u}, \mathbf{c}$

$$\sum_{e \in \{h, o\}} \int_0^T l^e(x, u, \mathbf{c}) dt$$



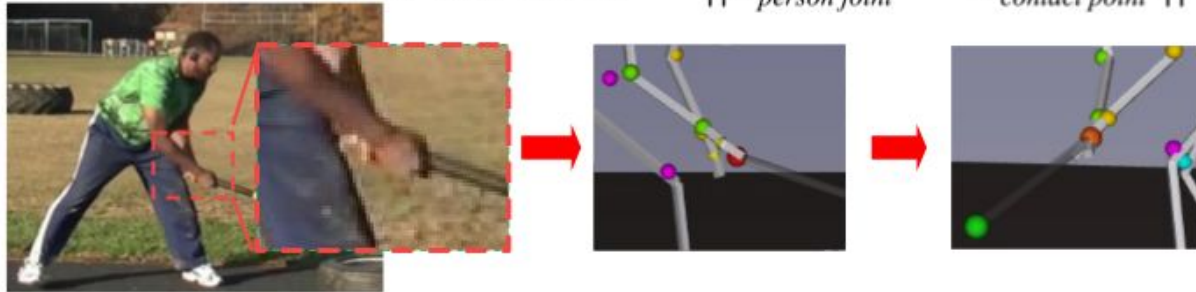
**Contact positions**

# Estimation Stage

$$\underset{\underline{x}, \underline{u}, \underline{c}}{\text{minimize}} \quad \sum_{e \in \{h, o\}} \int_0^T l^e(x, u, c) dt$$

Subject to

1. Contact motion model:  $\| p_{person\ joint} - p_{contact\ point} \| = 0$



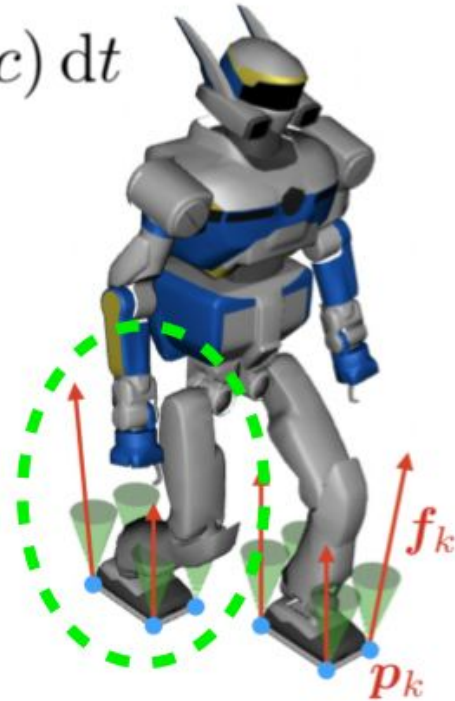
# Estimation Stage

$$\underset{\underline{x}, \underline{u}, \underline{c}}{\text{minimize}} \quad \sum_{e \in \{h, o\}} \int_0^T l^e(x, u, c) dt$$

Subject to:

2. Contact force constraints:

- Prevent the feet from sliding





# Estimation Stage

$$\underset{\underline{x}, \underline{u}, \underline{c}}{\text{minimize}} \quad \sum_{e \in \{h, o\}} \int_0^T l^e(x, u, c) dt$$

Subject to:

3. Lagrangian dynamics equation

$$M(\mathbf{x}) \ddot{\mathbf{x}} + \mathbf{b}(\mathbf{x}, \dot{\mathbf{x}}) = \boldsymbol{\tau}(\mathbf{u}, \mathbf{c})$$

[Carpentier, et al. Pinocchio. <https://stack-of-tasks.github.io/pinocchio>. 2015-2019]

# Composite Loss Function

- Loss function  $l^e$  = weighted sum of multiple costs
- **Data term** - observed vs re-projected 2D joint and obj endpoint positions.
- **Pose prior term** - encourage only physically plausible poses.
- **Physical plausibility term** - contact motions, contact forces, full body dynamics.
- **Trajectory smoothness term**

# Data Term

- Minimise the re-projection error of the estimated 3D human joints and 3D objection endpoints w.r.t 2D measurement in the video frame.
- Given  $\mathbf{j} = \mathbf{1}, \dots, \mathbf{N}$  joints or object endpoints

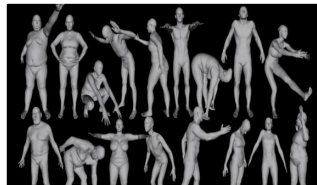
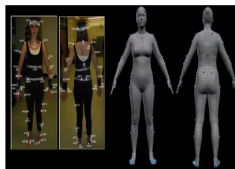
$$l_{data} = \sum_j \rho \left( p_j^{2D} - P_{cam}(p_j(q)) \right)$$

# Pose Prior Term

- **Motivation** - Single 2D skeleton = projection of multiple 3D poses; some unnatural or impossible exceeding human joint limits.
  - **Solution** - Limit to only plausible pose priors.
-

# Prior Pose Term (Obtaining pose priors)

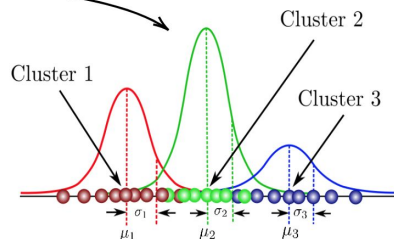
$MoSh =$



$= SMPL(MoSh($



$))$   
CMU MoCap



CMU Mocap data, <http://mocap.cs.cmu.edu/>

Loper et. al "MoSh: Motion and Shape Capture from Sparse Markers" SIGGRAPH Asia 2014"

Bogo et. al, "Keep it SMPL: Automatic Estimation of 3D Human Pose and Shape from a Single Image, ECCV 2016"

## Pose Prior Term

- Map  $q^h$  to SMPL pose vector and compute likelihood under pre-trained GMM
  - Minimise to favor more plausible poses vs impossible ones
- 

$$l_{pose}^h = -\log(p(q^h; GMM))$$

# Physical Plausibility Term

- **Contact motions** - minimise distance between active joints and contact points on object.
- **Contact forces** - model constraints induced by physical forces (friction, Newton's 3rd Law).
- **Full body dynamics** - model trajectory of body-object system using Lagrangian mechanics.



# Contact Motions

- For an active joint minimise distance between joint and active contact point

$$\|p_j^h(q^h) - p_k^c(x, c)\| = 0$$

- For planar contact motion

$$\kappa(x, c) = \sum_j \sum_{k \in \phi(j)} \delta_j \left\| T^{(kj)} (p_j^h(q^h)) - p_k^c(x, c) \right\|$$

# Contact Forces

- Environment exerts a contact force  $f_k$  on each active contact point.
- Two distinct contact forces.
- 6D spatial forces exerted by objects - 3D linear force and 3D moment.
- 3D linear forces due to friction - constrained inside 3D “cone of friction” approximated as 3D pyramid.

$$f_k = \sum_{n=1}^N \lambda_{kn} g_n^{(3)} \quad f_j = \sum_{k=1}^4 \begin{pmatrix} f_k \\ p_k \times f_k \end{pmatrix} = \sum_{k=1}^4 \sum_{n=1}^N \lambda_{jkn} g_{kn}^{(6)}$$

# Full Body Dynamics

- Fully body movements of person and manipulated object described by Lagrangian dynamics equation.

$$M(q)\ddot{q} + b(q, \dot{q}) = g(q) + \tau$$

# Trajectory Smoothness Term

- Regularise human and object motion.
- Leverage temporal continuity in the video to improve smoothness of human and object motion.
- Minimise spatial velocities and accelerations for smoothness and removing incorrect 2D poses.

$$l_{\text{smooth}} = \sum_j \left( \|\nu_j(q, \dot{q})\|^2 + \|\alpha_j(q, \dot{q}, \ddot{q})\|^2 \right)$$

- Both linear and angular movements are smoothed simultaneously.

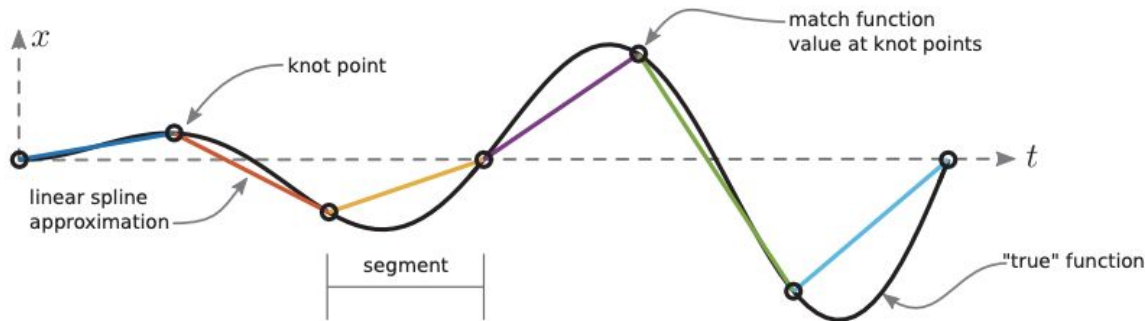
# Trajectory Smoothness Term

- Regularise contact motion and forces.
- Minimise the velocity of the contact points and temporal variation of contact force.

$$l_{\text{smooth}}^c = \sum_j \sum_{k \in \phi(j)} \delta_j \left( \omega_k \|\dot{c}_k\|^2 + \gamma_k \|\dot{f}_k\|^2 \right) dt$$

# Optimisation

- Convert continuous problem to discrete non-linear optimisation problem using collocation method.
- Discretise all trajectories -  $x$   $c$   $u$



# Extracting 2D Measurements from Video

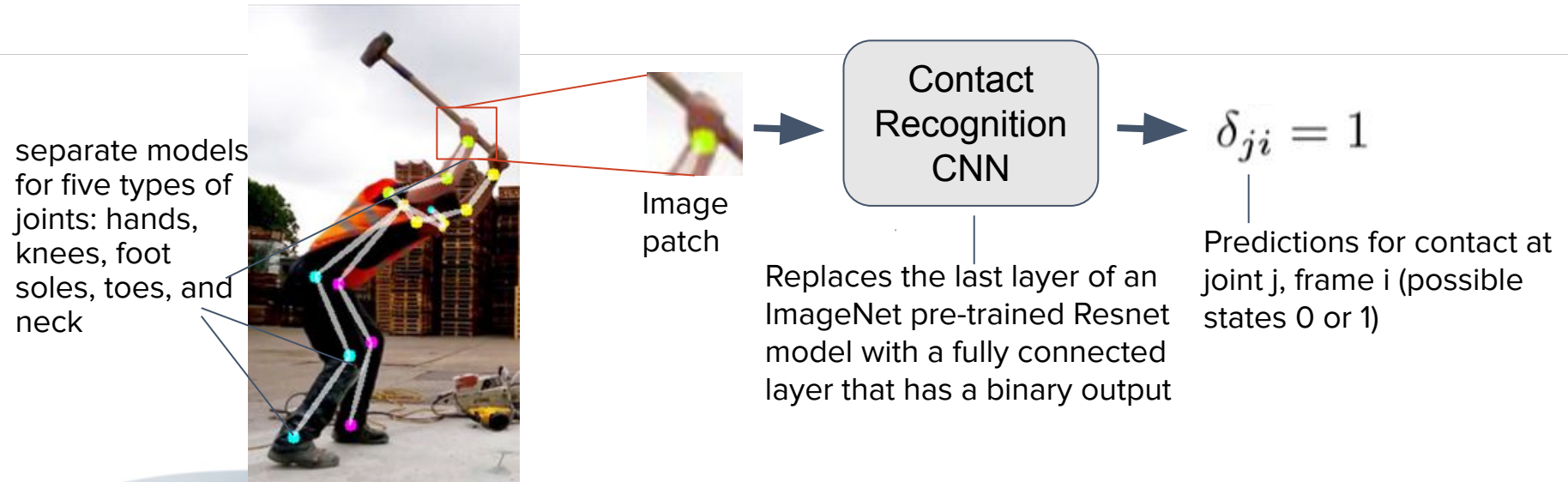
- **Estimating 2D position of joints:** OpenPose [Zhe Cao et al., Berkeley AI Research]





# Extracting 2D Measurements: Contact Points

- Training dataset: manually annotated contact data (Google Image Search & youtube videos) harvested from the Internet



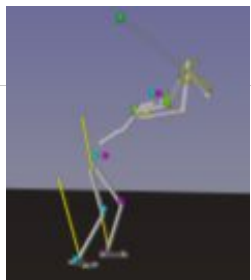


---

# Experimental Results

# Parkour Dataset

## Mean per Joint Position Error (in mm)



Method	Jump	Move-up	Pull-up	Hop	Avg
SMPLify [9]	121.75	147.41	120.48	169.36	139.69
HMR [37]	111.36	140.16	132.44	149.64	135.65
Ours	<b>98.42</b>	<b>125.21</b>	<b>119.92</b>	<b>138.45</b>	<b>122.11</b>

Improvement over state-of-the-art for Barbell, Spade, and on Average over all tool types.

# Parkour Dataset



## Errors of the contact forces for soles and hands

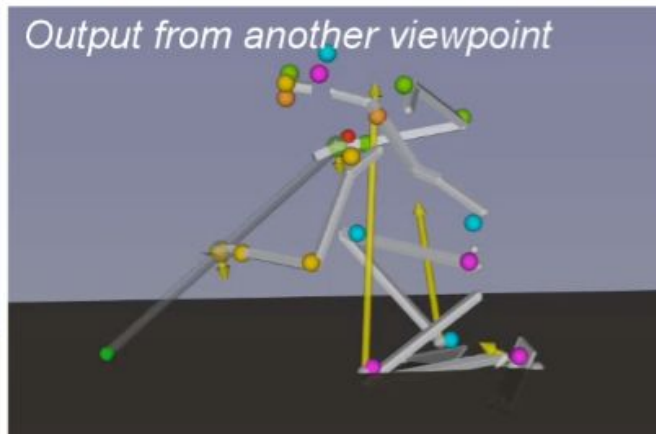
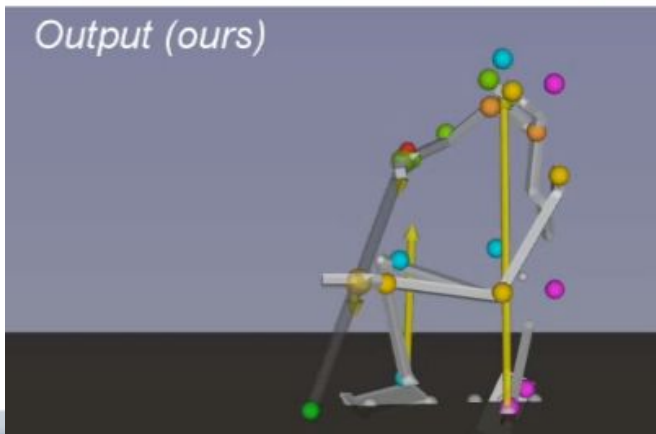
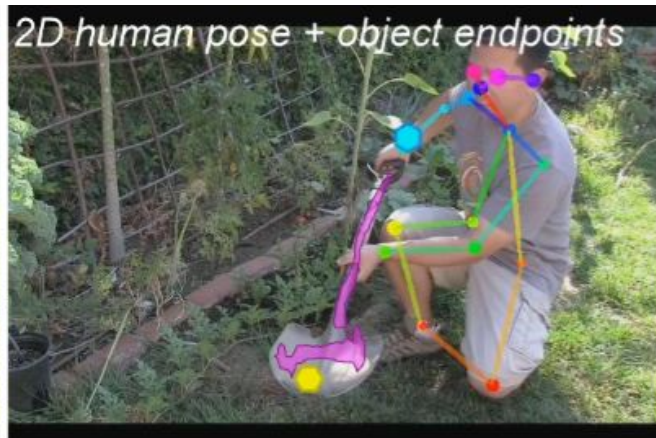
	L. Sole	R. Sole	L. Hand	R. Hand
Force (N)	144.23	138.21	107.91	113.42
Moment (N·m)	23.71	22.32	131.13	134.21

To estimate the forces we assume a generic human physical model of mass 74.6kg for all the subjects

# A New Dataset: Handtool Videos







# Results: Video

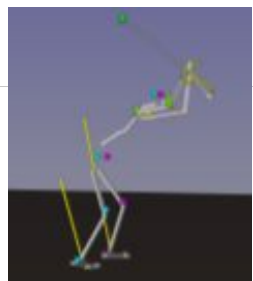
[Video](#)

---



# Results: 3D Pose Estimation on Handtool Dataset

## Mean per Joint Position Error (in mm)



Method	Barbell	Spade	Hammer	Scythe	Avg
<b>SMPLify</b>	130.69	135.03	93.43	112.93	118.02
<b>HMR</b>	105.04	97.18	96.34	115.42	103.49
<b>Ours</b>	104.23	95.21	95.87	114.22	102.02

Improvement over state-of-the-art for Barbell, Spade, and on Average over all tool types.

# Results: 2D Endpoint Estimation

% of Endpoints within 25/50/100 pixels from ground truth



<b>Method</b>	<b>Barbell</b>	<b>Spade</b>	<b>Hammer</b>	<b>Scythe</b>
<b>Mask-R-CNN</b>	33/42/54	54/79/93	35/44/45	63/72/76
<b>Ours</b>	38/71/98	57/86/99	61/91/99	69/88/98

Improvement over state-of-the-art for all tool types.

# Strengths & Weaknesses

- + Inducing physics priors to trajectory estimation.
- + Recovers both human-object interaction and forces that give rise to the motion.

---

- + Able to generalise to the Parkour dataset videos of 400-2200Hz although trained on  $\sim 25$ Hz videos.
- - More visual illustrations for the physical constraints would have helped.
- - Formulation of 6D force generators is difficult to understand (even in the Appendix).

# How can this work be extended?

- Articulating tools.
- Apply to robots.
- Marker-less motion capture data from videos
- What other data-driven/DL problems can benefit for infusion of common sense priors (e.g. physics, optics) as part of the optimisation stage?

---

TGIF!!!)