# Sign language coding, 3D behavior data ... and ANVIL

**Michael Kipp**

DFKI
Embodied Agents Research Group
Cluster of Excellence
*Multimodal Computing and Interaction*
Saarland University

**OTIM / ILIKS Workshop**
24 May 2011
LPL, Aix-en-Provence

*Joint work with:*
Alexis Heloir, DFKI
Quan Nguyen, DFKI
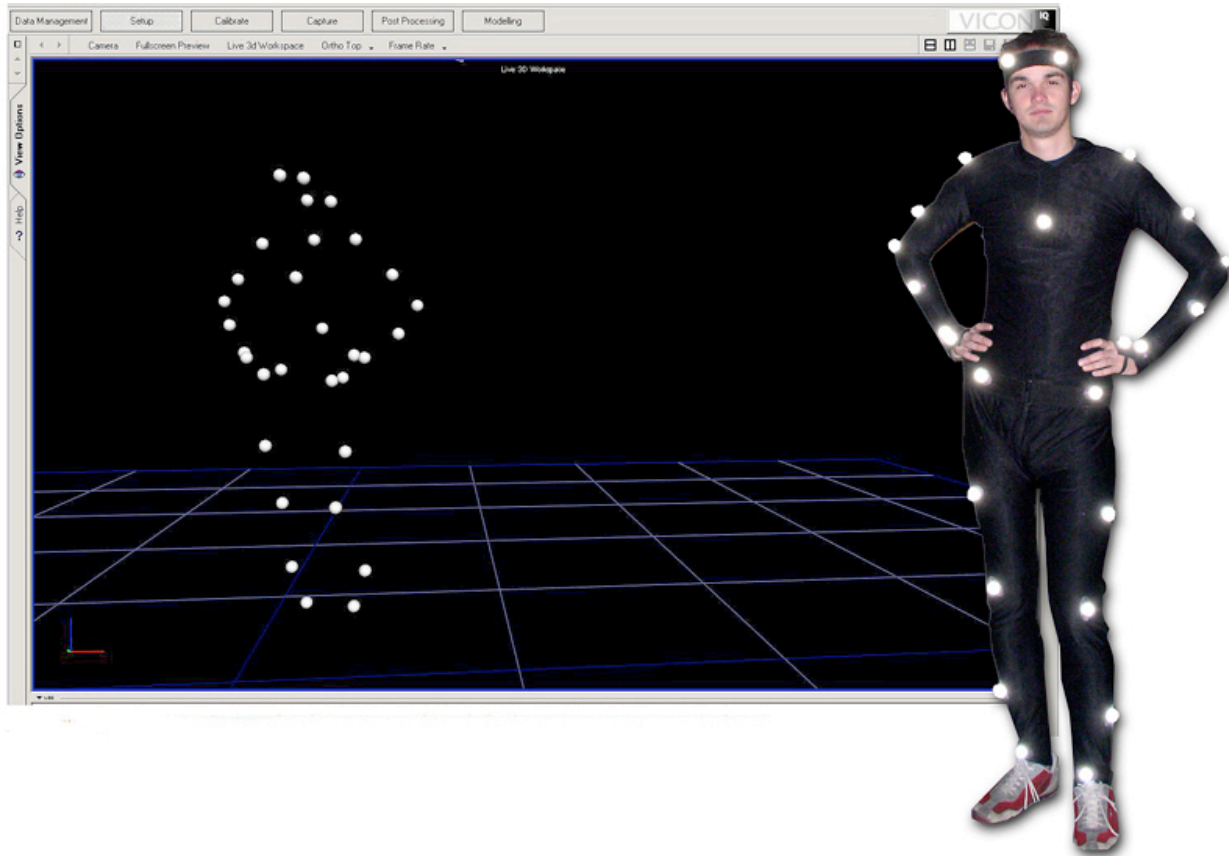Michael Neff, UC Davis

# Overview

- Multimodal corpora for animation

- Sign language avatars

- ANVIL

**Announcements:**
Workshop on Multimocal Corpora: Taking Stock and Roadmapping the Future
held in conjunction with ICMI-2011 (Heylen, Paggio, Kipp), 18 November
Watch www.multimodal-corpora.org

Second International Workshop on Sign Language Translation and Avatar Technology (SLTAT),
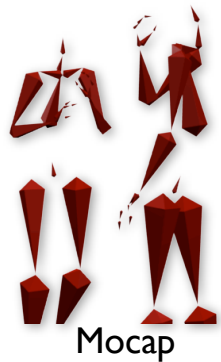with ACM ASSET 2011, Dundee, UK.
Watch http://embots.dfki.de/SLTAT

# Corpora for Animation

# What Can be Learned from Motion Data ?
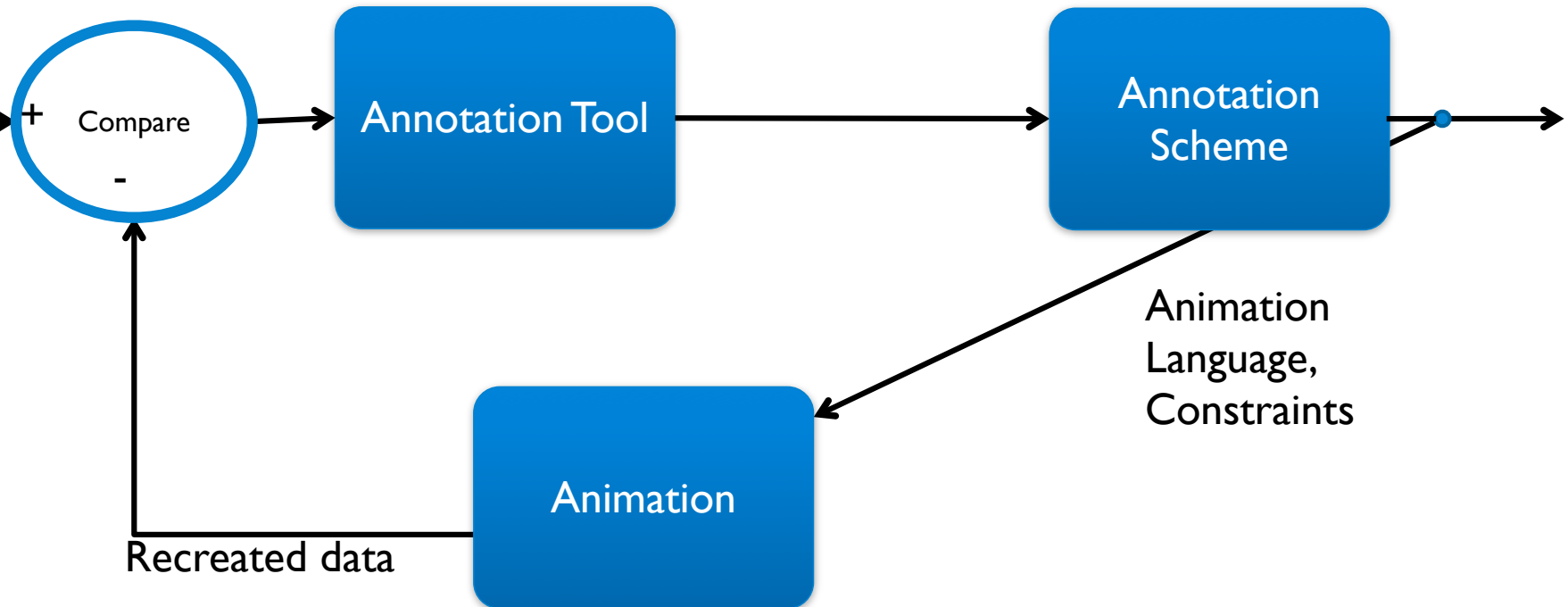
▸ Ambient movements (Egges et al. 2005)

▸ Balance control (Neff et al. 2009)

▸ Motion graphs (Kovar et al. 2002)

▸ Recreation of gesture from annotations (Kipp et al. 2008)

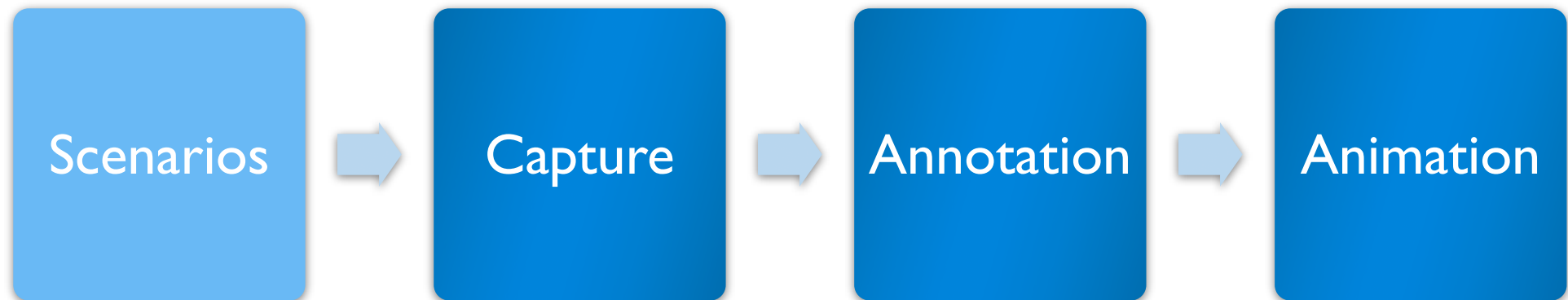▸ We are interested in building generative models of communicative gestures (in dyadic conversations)

▸

# The analysis and synthesis loop

**How can empirical data improve animation methods?**

Mocap

Video

Compare

+

-

Annotation Tool

Annotation Scheme

Animation

Recreated data

Animation Language, Constraints

# Data Acquisition and Processing Pipeline

Scenarios → Capture → Annotation → Animation

# Video Corpus
# (Neff et al. 2008)



# Mocap/Video Corpus
# (Heloir et al. 2010)

# Recent Capture Session (UC Davis)

- Improvised acting
- 19 dyadic scenarios (two friends meet …)
- Status high/low + agree/disagree
- Proxemic behavior + NVB synchronization in dyads
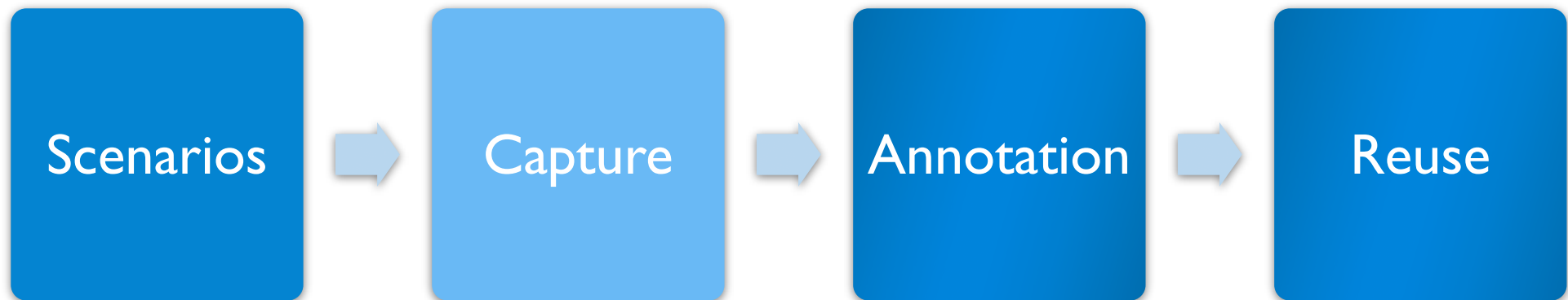


Status



Liking

# Example: Two People Meet
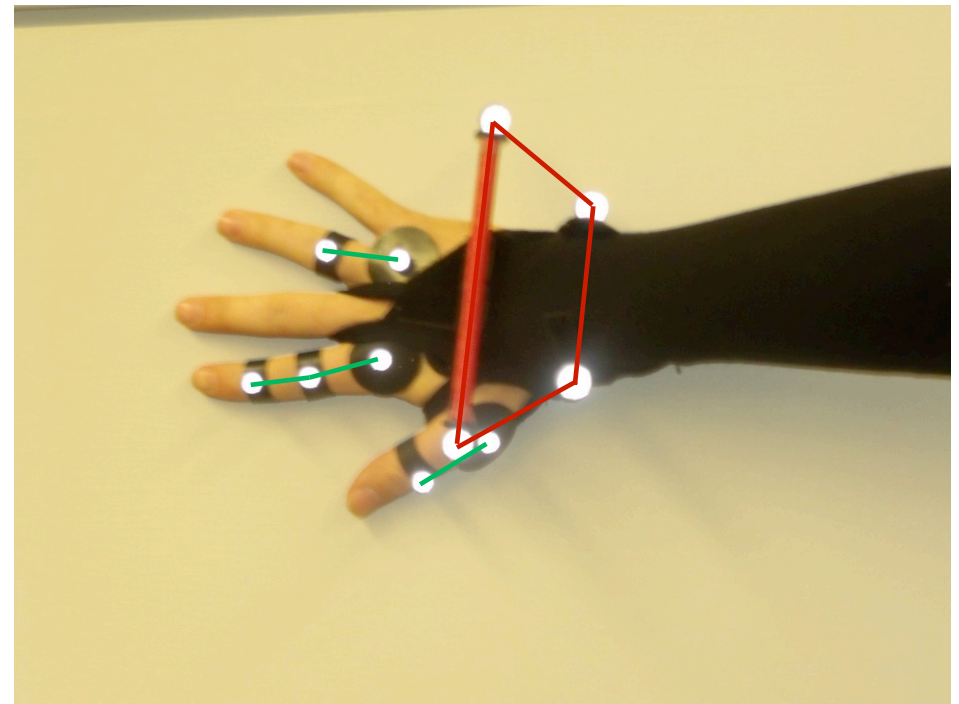


They like each other

They dislike each other

# Technical setup

- Optical Motion Capture
  - Vicon MX 40
  - 12 Cameras

- Video recorder (x2)
  - HD

- Camera mounted
  microphones
  (not recommended)
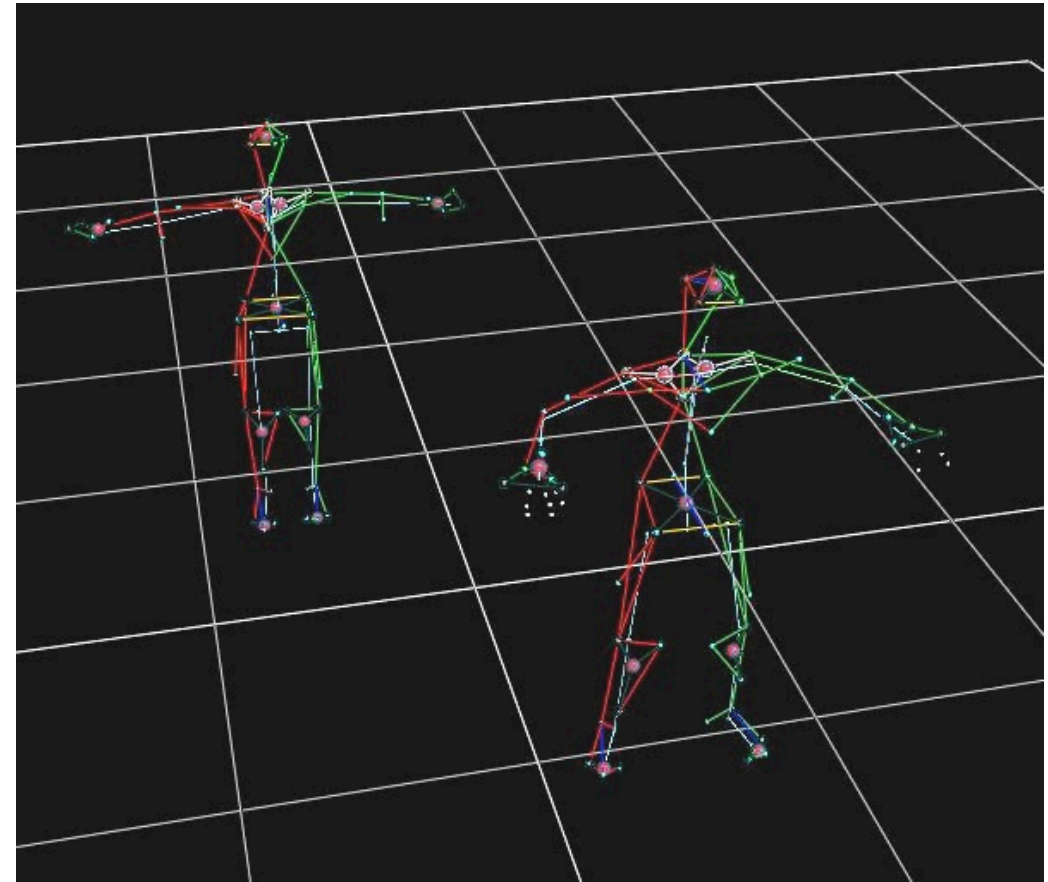
# Capturing Handshape

▸ Occlusions are frequent between fingers

▸ Impossible to record motion for all fingers

▸ We used a reduced set of markers
  ▸ index finger
  ▸ thumb
  ▸ „rest"



▸ **Similar to** (Chang et al., 2007)

▸

# Reconstruction

- From marker clouds to skeleton

- Semi automated process

- Significant manual processing required
  - labeling correction (occlusion, confusion, mainly hands)
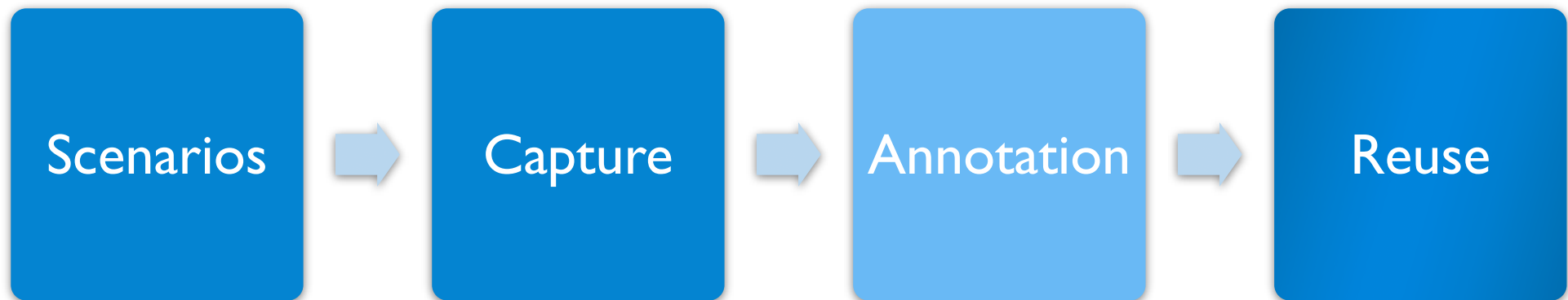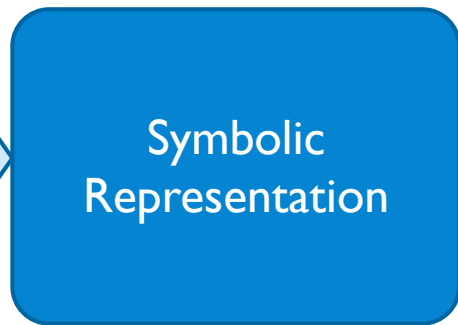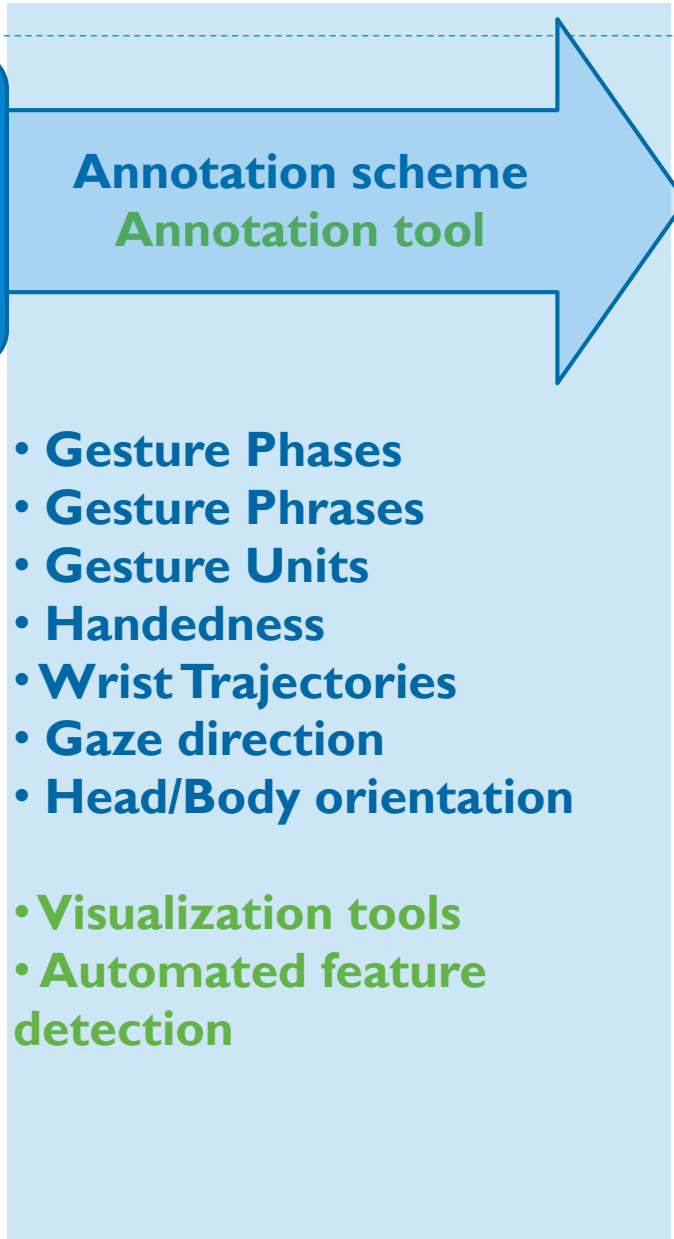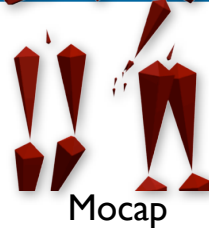
- Postprocessing work: 1 : 40

for instance: BVH files

# Annotation: From Raw Data to Symbolic Representation

**Raw Data**

Mocap

Video

- fine-grained
- no „meaning"
- difficult to manipulate
- highly realistic

**Annotation scheme**
**Annotation tool**

- **Gesture Phases**
- **Gesture Phrases**
- **Gesture Units**
- **Handedness**
- **Wrist Trajectories**
- **Gaze direction**
- **Head/Body orientation**

- **Visualization tools**
- **Automated feature detection**

Symbolic Representation

```
BEGIN K_POSE_SEQUENCE
  CHARACTER:Amber
  START:asap
  FADE_IN:200
  FADE_OUT:200
  BEGIN K_POSE          # first pose lower middle
    TIME_POINT:+800
    HOLD:200
    BEGIN POSITION_CONSTRAINT
      BODY_GROUP:rarm
      TARGET:-0.07;-0.3;0.13
      JOINT:rhand
      OFFSET:0.0;0.0;0.0
    END
    BEGIN ORIENTATION_CONSTRAINT
      BODY_GROUP:rarm
      NORMAL:Zaxis
      DIRECTION:0.0;0.0;-1.0
      JOINT:rhand
    END
    BEGIN POSITION_CONSTRAINT
      BODY_GROUP:larm
      TARGET:0.07;-0.3;0.13
      JOINT:lhand
      OFFSET:0.0;0.0;0.0
    END
    BEGIN ORIENTATION_CONSTRAINT
      BODY_GROUP:larm
      NORMAL:Zaxis
      DIRECTION:0.0;0.0;-1.0
      JOINT:lhand
    END
    BEGIN SWIVEL_CONSTRAINT
      BODY_GROUP:rarm
      TARGET:30
      JOINT:rhumerus
    END
    BEGIN SWIVEL_CONSTRAINT
      BODY_GROUP:larm
      TARGET:30
      JOINT:lhumerus
    END
    BEGIN ORIENTATION_CONSTRAINT
      BODY_GROUP:larm
```
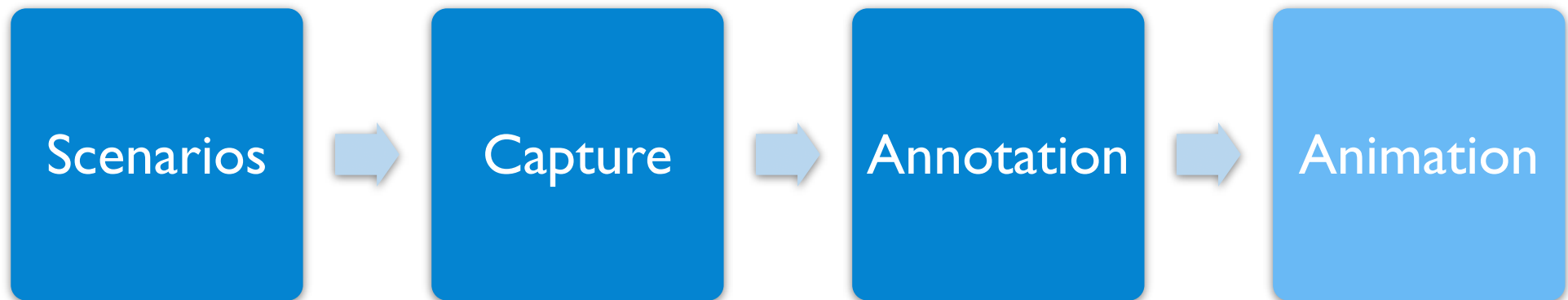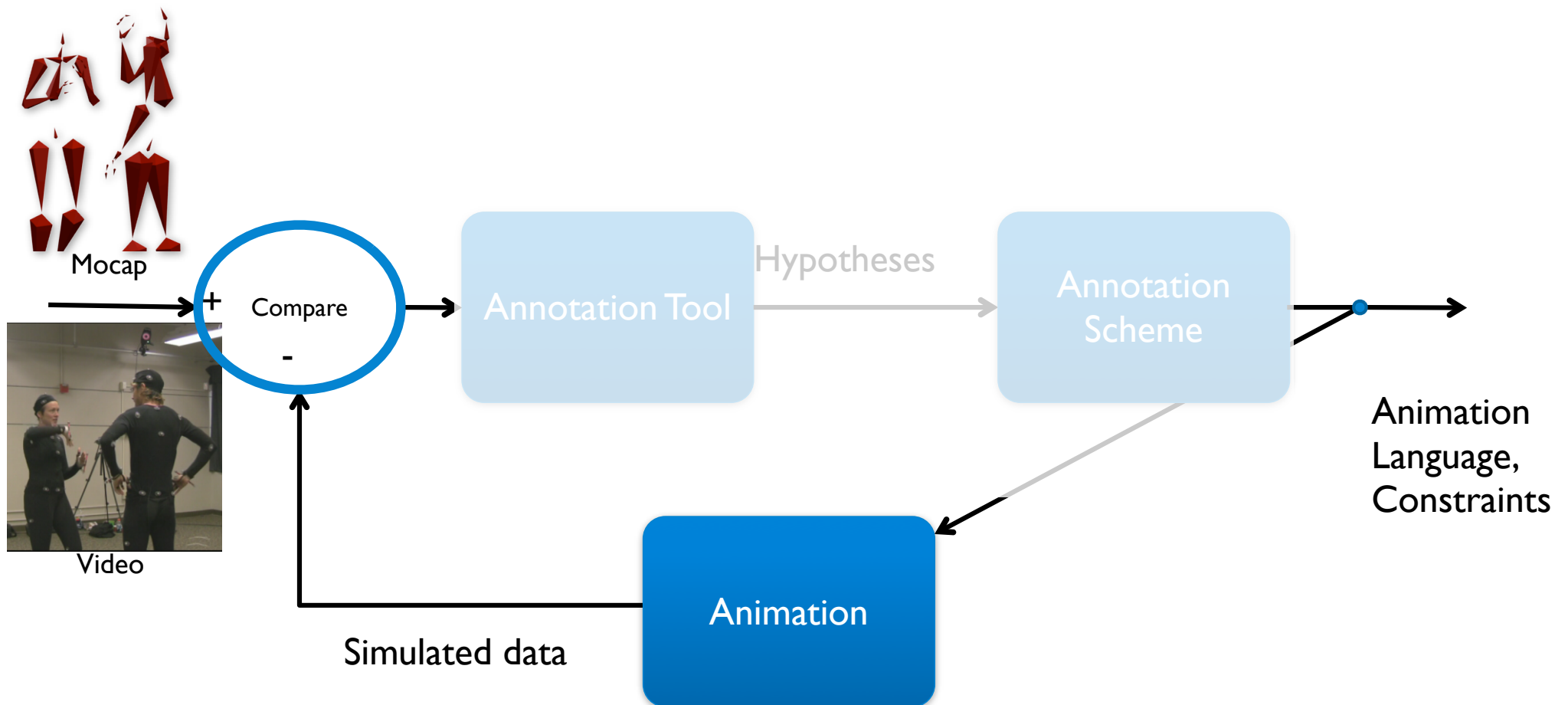
Animation Language, set of constraints

- compact
- meaningful
- easy to manipulate
- realistic ???

# Validation by Recreation



Mocap

Video

Compare + −

Annotation Tool

Hypotheses

Annotation Scheme

Animation Language, Constraints

Animation

Simulated data
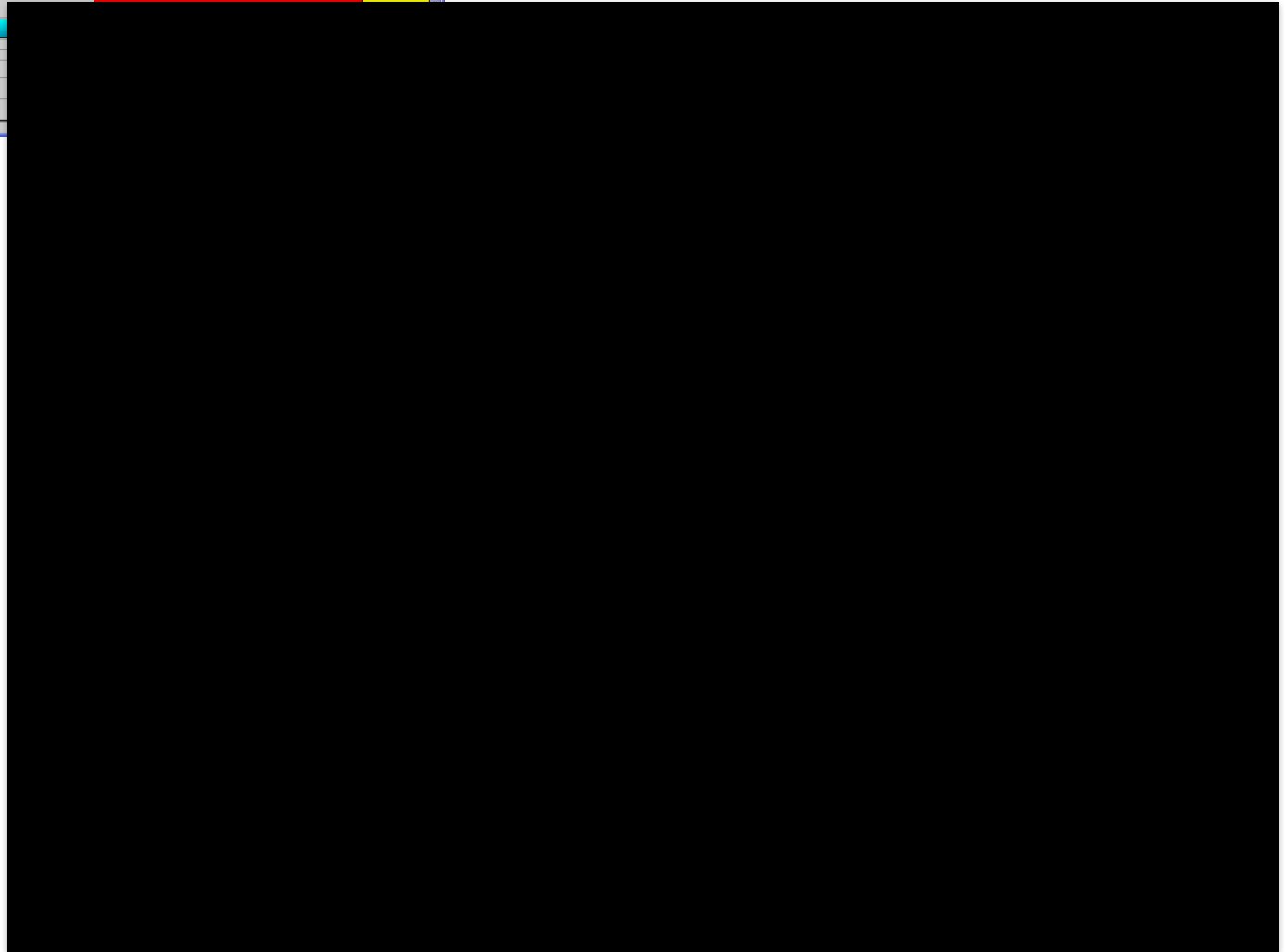
Just „recreated"
(some call it „reanimated")

# Gesture Modeling and Animation Based on a Probabilistic Recreation of Speaker Style

[Neff et al. 2008] ACM Transactions on Graphics
[Kipp et al. 2007] JLRE (coding scheme)

- Question: Longer **G-Units** => more natural?
- Hypothesis: Yes
- Experiment
  - **G version**: synthesized
  - **S version**: manipulation (made singular)



- Results
  - **G version**:
    - **more natural**            $p < .01$
    - **more friendly**            $p < .001$
    - **more trustworthy**      $p < .001$
  - **S version**:
    - **more nervous**           $p < .001$

[Kipp et al. 2007] IVA 07, Best paper

# Why Motion Capture?



Video

vs.

signal
=>

spectrogram
waveform
intensity
pitch
=>
segmentation
categorization

projection onto a 2D screen
merged with background
=> degraded, noisy signal
=> seg. + cat.

# Why Motion Capture?

- **Objective measures**

  ➡ speed / velocity (acceleration)

  ‣ rhythm analysis, interpersonal synchrony, correlation with intonation

  ➡ shape of the gesture

  ‣ trajectory, motion contour

  ➡ hand location in gesture space (automatic/robust)

  ➡ direction of a gesture (vector)

  ➡ distance and orientation of interlocutors (proxemics)

- **Viewing**

  ➡ Watch from any angle

  ➡ Zoom in/out without quality loss!

  ➡ Virtual world visualization support
  (motion trails, coordinate system, boundary planes, vector arrows)

- **Automation** (segmentation, categories)

# ANVIL for Gesture Annotation

▶ supports motion capture
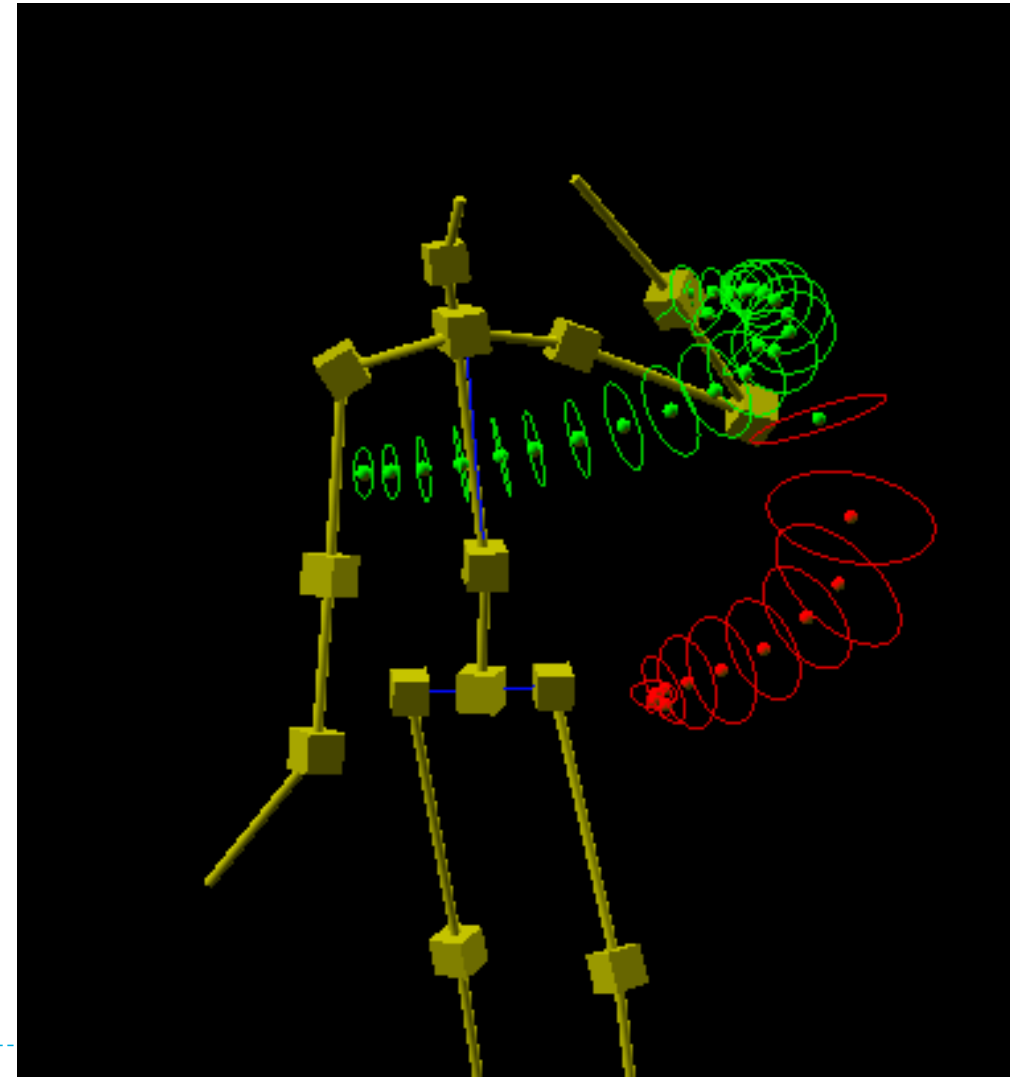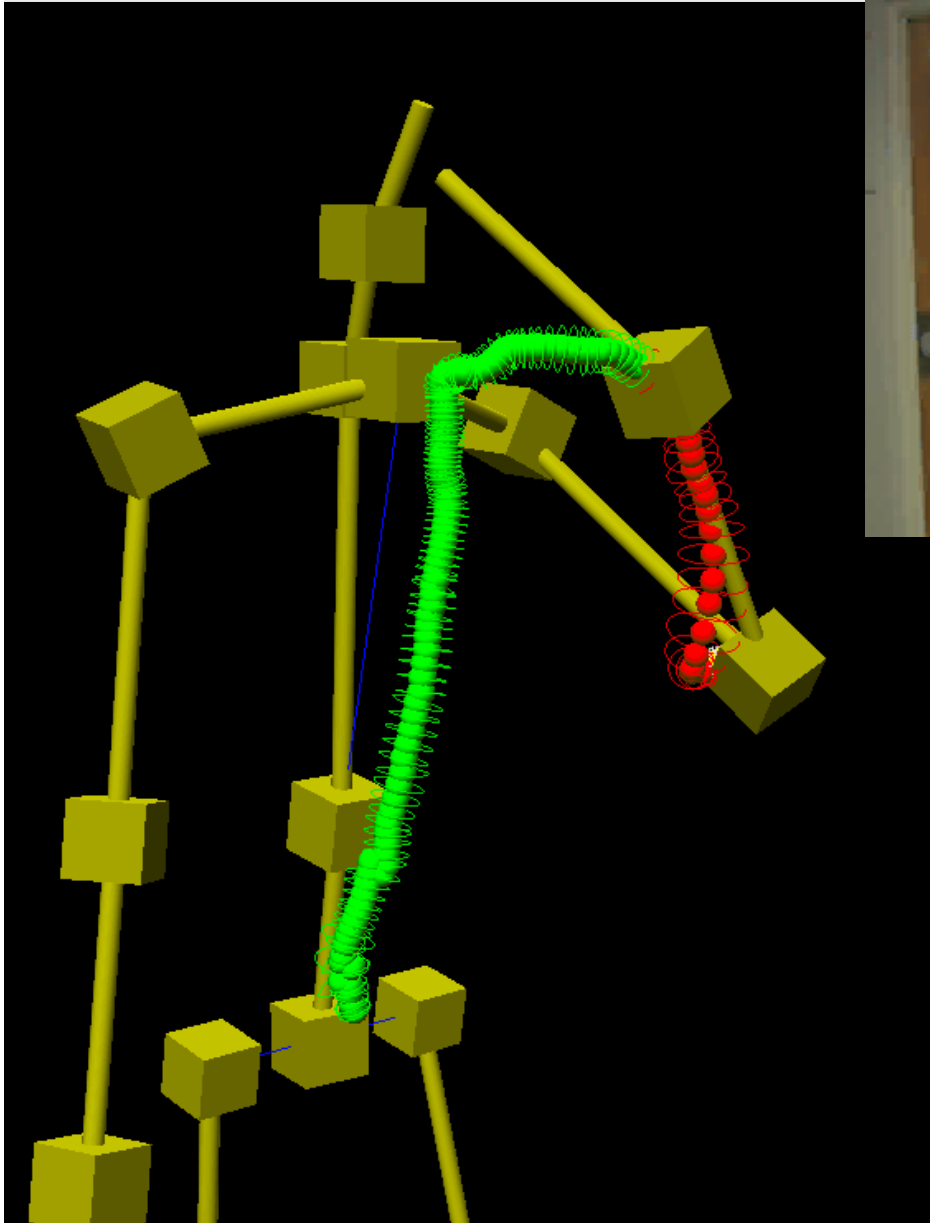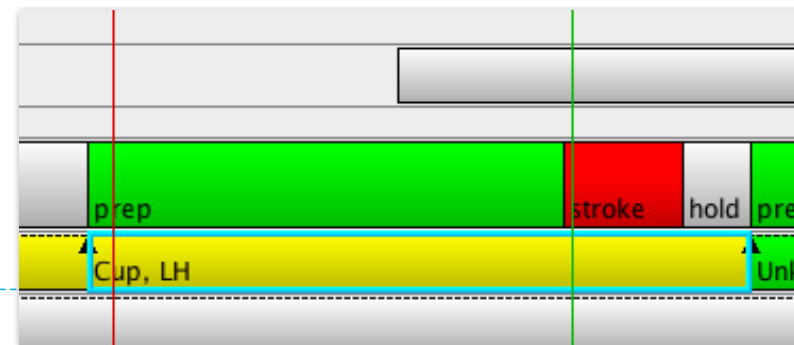
  ▶ synchronization of video, sound and mocap

# Motion Trails

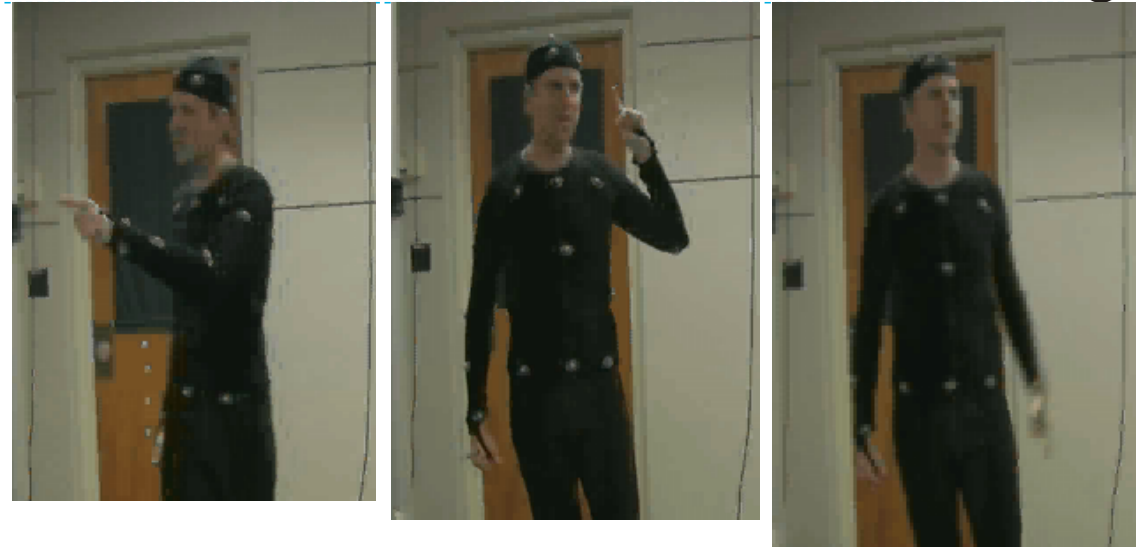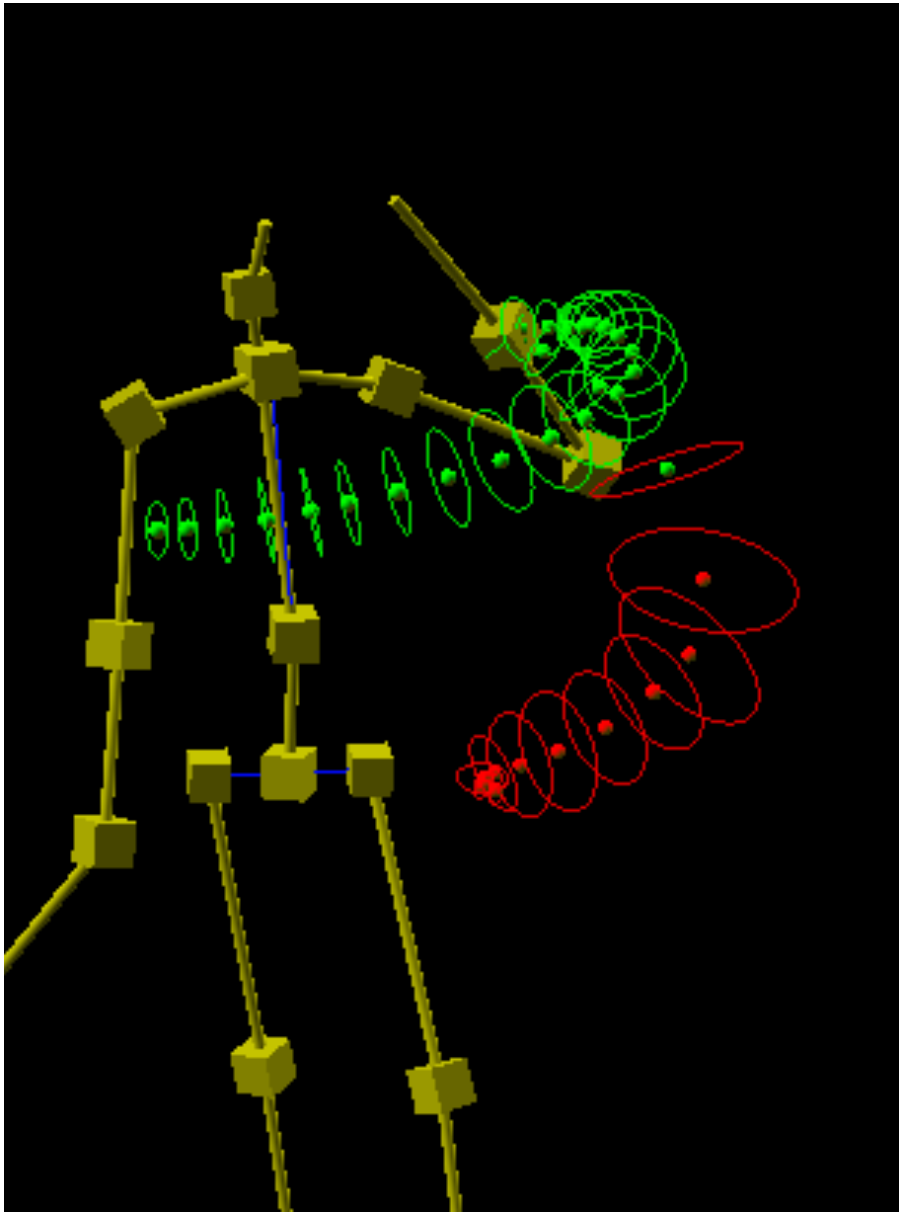- Continuous representation of motion in 4D

- Shows segmentation by color coding

- Gives an impression of the velocity profile
  - spacing: indicator but too subtle
  - circles: indicate direction vector, can be scaled (gain), do not occlude

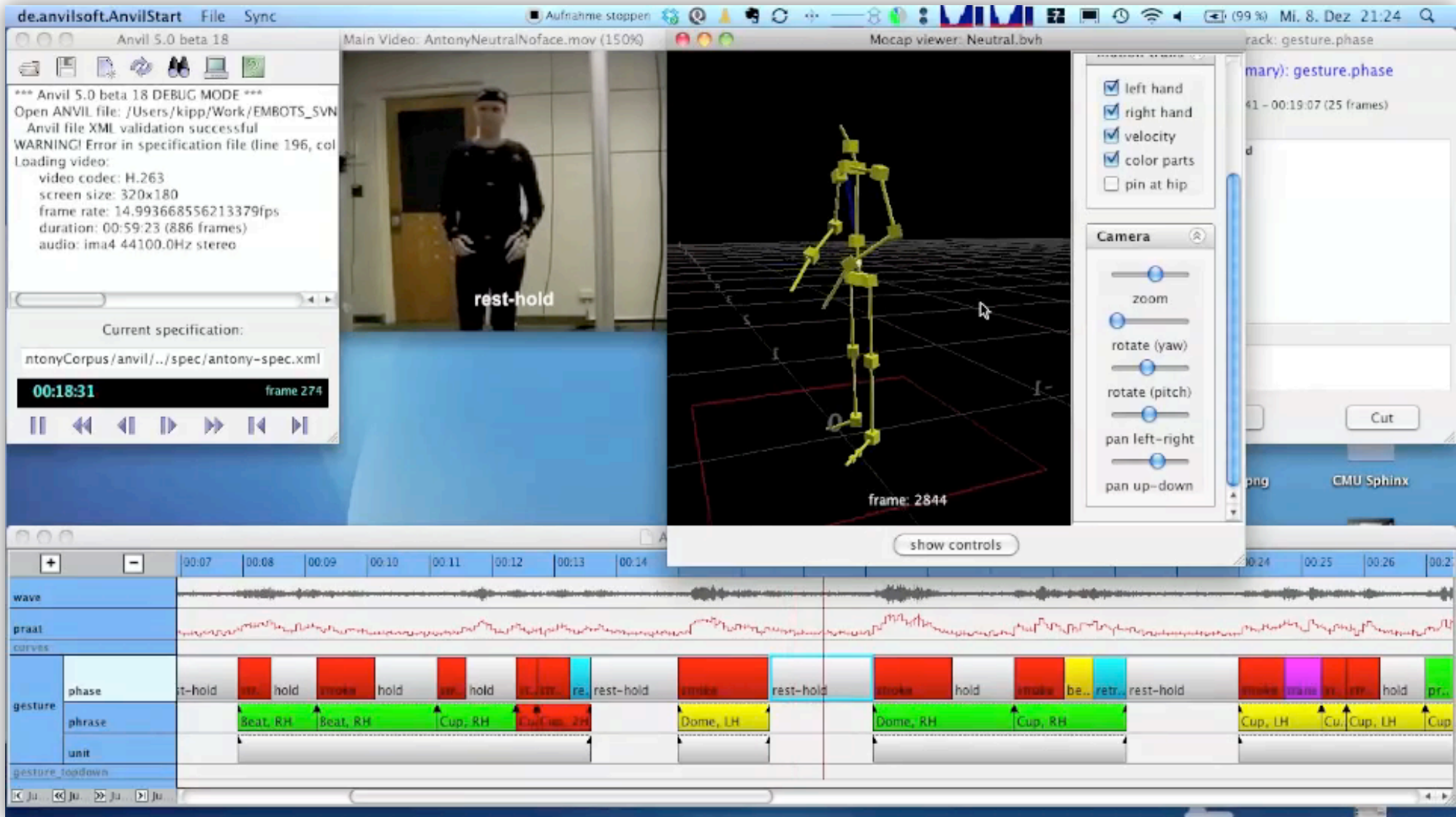- Prevents annotation errors
  - incorrect hand or omission

- Slow gesture → small circles, spheres close to each other
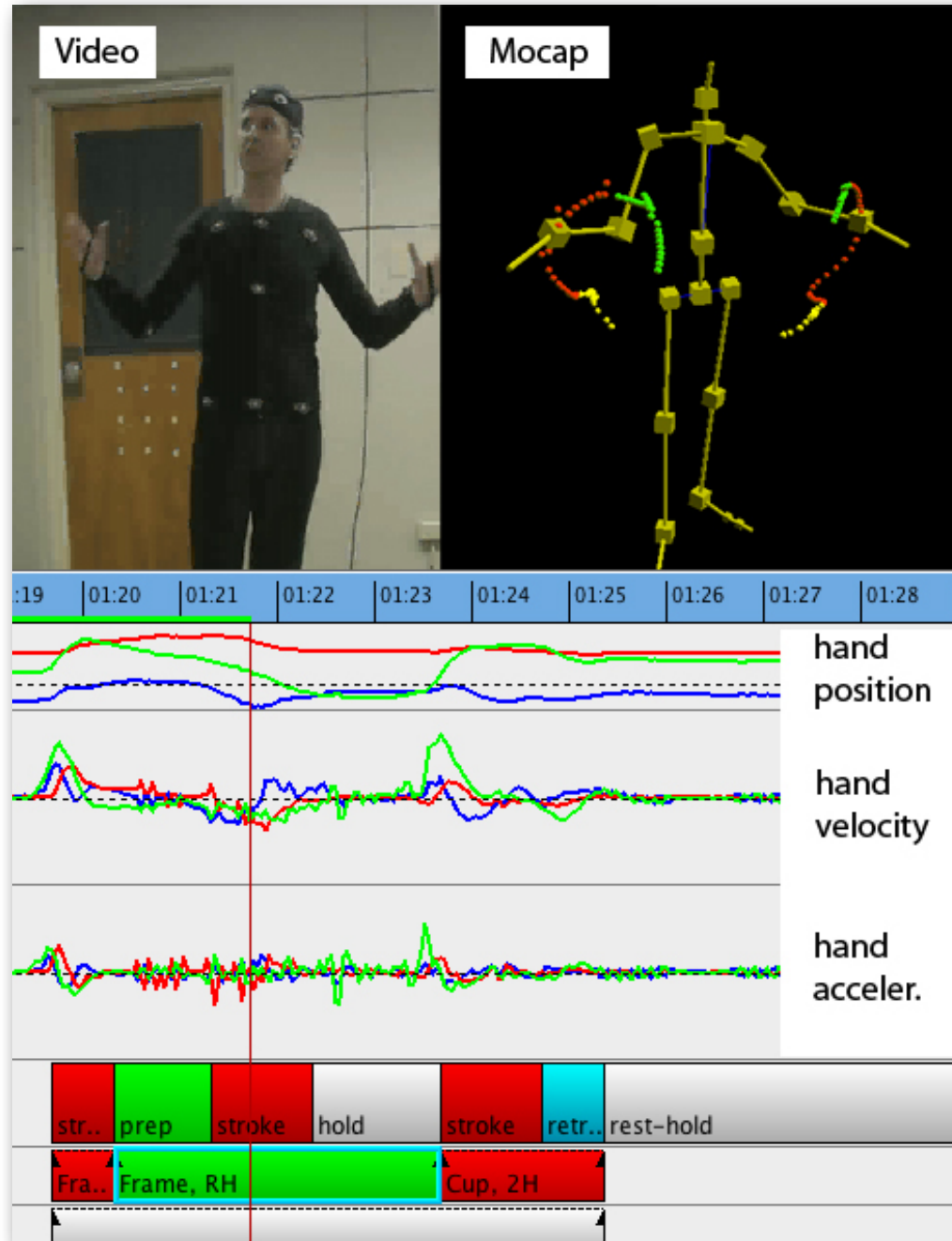- Trail shows sideways motion of hand

- High velocity → bigger circles, spheres more spaced
- Shows nicely that **stroke** has even higher velocity at ist peak

# Video + Mocap + 4D Trails



[Kipp 2011, Heloir, Neff, Kipp, 2010]

# Automated Handedness Detection

▸ Handedness automatically detected gesture-wise
   (= single gesture)

▸ Compare the normalized path length of right hand $L_{RH}$ and
   left hand $L_{LH}$ over a gesture of length d (time)
   → if $\dfrac{|L_{RH} - L_{LH}|}{d} < 0.12\dfrac{m}{s}$ then gesture labeled 2H

▸ First test on **269** gestures: 83% correct

So you don't have 100,000$ for your personal mocap lab?

▸

# Poor Man's Mocap

- Microsoft Kinect

- Hacked one hour after release

- Free software for skeleton tracking

- No excuses :)

**Recipe „Kinect for Anvil"**

1. Install various software
   (OpenNI, NITE, Brekel)
2. Plug in kinect
3. Calibrate in Psi pose
4. Switch on video camera
   (for later coding)
5. Click „start capture bvh"
6. Load everything in ANVIL
7. Sync video and mocap

124.90 €

XBOX 360

PrimeSense™
Natural Interaction

=> can produce a .bvh file (demo)
     we will put a „howto" online soon
=> view & annotate in ANVIL

**DEMO**

mydata.bvh     Biovision     **ANVIL**
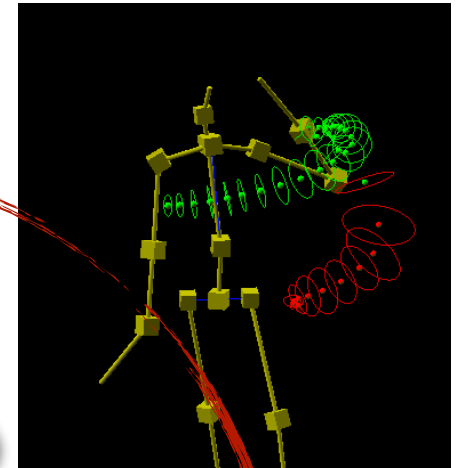
hierarchical skeleton definition

motion data

RHand

Hip

LHand

Joint structure

One line per „frame": angles for all joints

# Signing Avatars

# Sign Language Avatars

- 500,000 Deaf in Europe

  ➡ Mother language / primary language: sign language

  ➡ Spoken language = second language (hard to learn!)

  ➡ 80% of the deaf leave school with significant reading/writing problems

- German Federal Ministry for Labour and Social Affairs (BMAS):

  ➡ Are signing avatars a solution for accessible dynamic web content? (current comprehensibility around 60%)

- Feasibility study

  ➡ state of the art, research priorities, applications

# Focus Group Interviews

# Existing Avatars

# Avatar videos: criticism

- Upper body

  ➡ too little involvement, especially no sideways rotations

  ➡ important in constructed dialogue

- Face

  ➡ too little eyebrow movement

  ➡ hardly any mouthing (important for DGS)

    ‣ absence of lip movement more striking than bad lip movement

    ‣ recent CG movie („Lissy") allowed quite good lip reading!

  ➡ missing teeth and tongue
  (necessary e.g. for letters L and N)

# Avatar videos: criticism

- **Style**

  ➡ hardly any emotional expression

  ➡ stiff / robotic movements

  ➡ missing personality easily interpreted as cold / unfriendly

- **Synchronization**

  ➡ mouthing and signs durations did not match

  ➡ important for keeping the face as a focus point (!), otherwise focus occillates between hands and face

# Avatar videos: criticism

- **Technical**

  ➡ good lighting and contrast important (e.g. black clothes are good)

  ➡ shadows support 3D effect
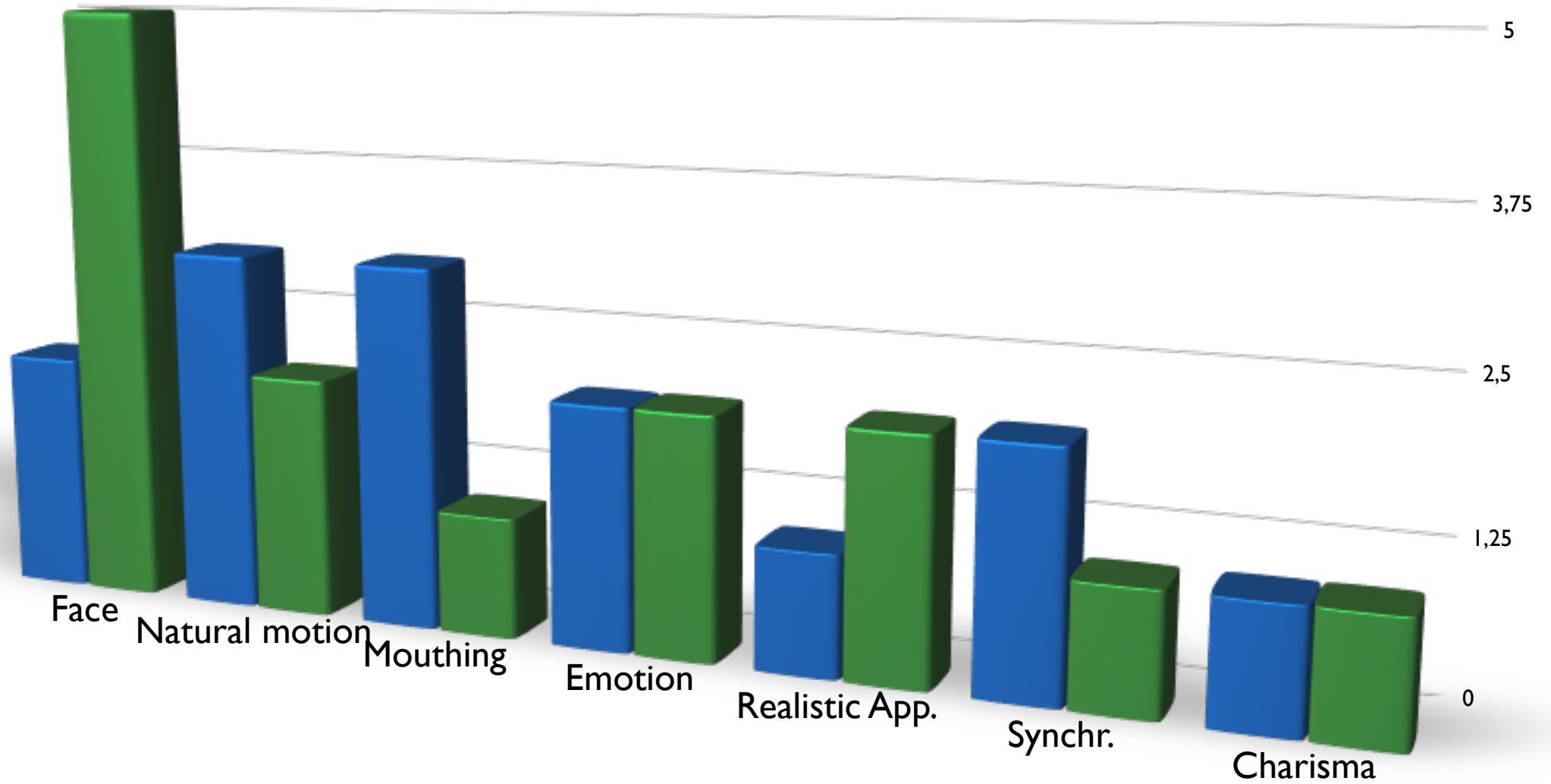
  ➡ preferably: speed & perspective under user control

- **Avatar Appearance**

  ➡ different avatars for different domains

  ➡ child avatar & cartoonish: for kids and entertainment

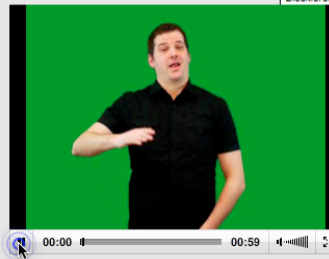  ➡ adult & realistic:  serious applications (politics, church ...)

# Avatar Aspects

Group 1 ■ Group 2

5

3,75

2,5

1,25

0

Face · Natural motion · Mouthing · Emotion · Realistic App. · Synchr. · Charisma

## Herzlich Willkommen



Wir möchten Sie ganz herzlich zu unserer Befragung begrüßen!
In dieser Befragung möchten wir Ihre Meinung und Einschätzungen zum Einsatz von Avataren als Gebärdensprachavatar wissen. Wir zeigen Ihnen dazu Videos und Bilder von Avataren, die Sie dann nach verschiedenen Kriterien bewerten sollst. Daher ist kein Vorwissen notwendig.
Im Abschnitt unten erklären wir Ihnen in der Projektbeschreibung, was wir machen und wozu wir diese Umfrage durchführen.
Die Befragung dauert ca. 20 min. Wir würden uns freuen, wenn Sie Ihre Antworten eventuell mit kurzen Sätzen oder Stichworten begründen könntest.

Wir bedanken uns ganz herzlich für Ihre Zeit und wünschen Ihnen viel Spaß bei der Befragung.

## Projektbeschreibung



Ein Avatar ist eine künstliche Figur in einer virtuellen Welt. Avatare könnten eingesetzt werden, um dynamische Texte von Internetseiten automatisch in Gebärdensprache übersetzen zu lassen. Dies könnte eine erfolgversprechende, langfristige Lösung sein, um Internetangebote für gehörlose Menschen barrierefrei zu gestalten. Bisher liegt die Verständlichkeit von Avataren jedoch nur bei ca. 60%. Bei einer verbesserten Verständlichkeit könnten die Einsatzmöglichkeiten von Gebärdenavataren zusätzlich ausgeweitet werden. Weitere mögliche Einsatzgebiete von Avataren könnten sein:
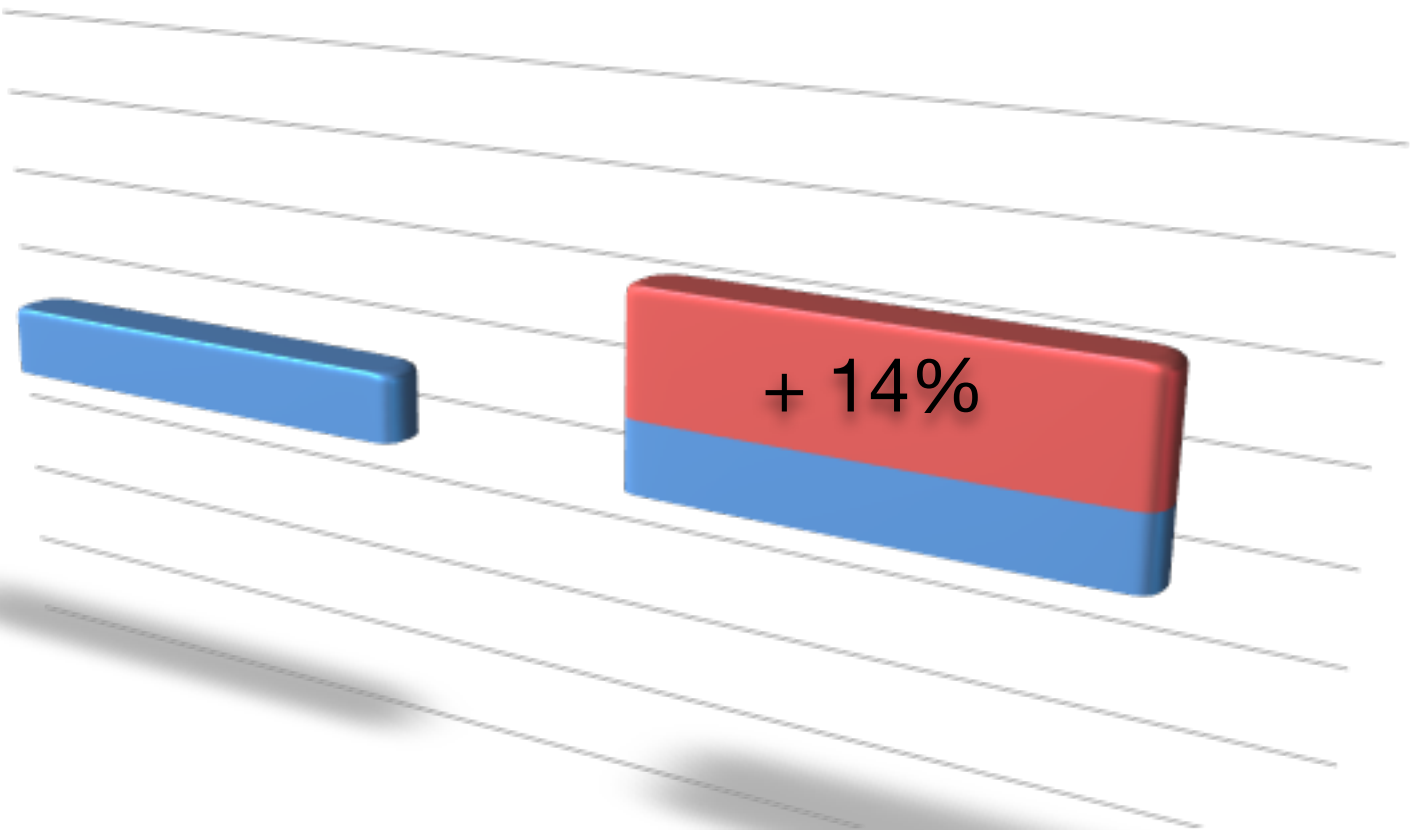1. Helfer bei Alltagssituationen (wie Zahnarztbesuch)
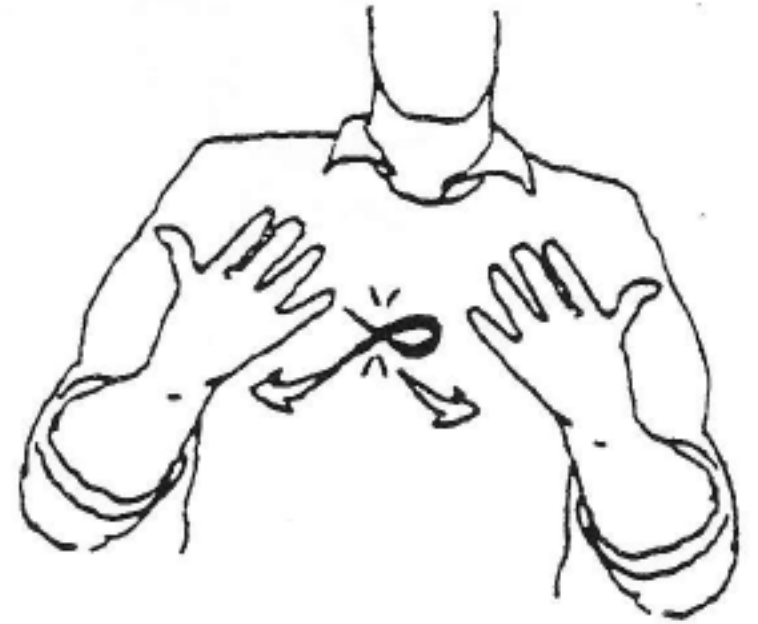2. Jobsuche
3. Wohnungssuche

In unserer Machbarkeitsstudie möchten wir eine kritische Bestandsaufnahme machen und mögliche technische Entwicklungen zusammenfassen. Dadurch sollen die Möglichkeiten und Grenzen des Einsatzes von Gebärdenavataren besser

# Do you consider avatars useful?
(-2 ... +2)

- Before: +0.4

- After: +0.7

+ 14%

Avatar
Sign Language
Animation

# Character Animation

- Initial Motivation

  ➡ Create a reusable character animation engine

  ➡ Exploration of coverbal gesture

- Existing systems (e.g. Greta, SmartBody, MAX) proposed high-level control languages

  ➡ Behavior Markup Language (BML)

  ➡ e.g. <gesture type="pointing" stroke="x" />

- Useful layer of abstraction but...

  ➡ What if you need more control (hand shape, torso involvement, size of gesture...) ?

# EMBR:
# [EM](EM)bodied Agent [B](B)ehavior [R](R)ealizer

- Our solution: Low-level control language

  ➡ Wrapper around keyframe animation

  ➡ Theory-independent (bottom-up approach)

- Used to create gesture lexemes

  ➡ Add „semantics" inside the gesture,
  e.g. this is the stroke, this is a preparatory motion

  ➡ Also: specify open parameters like hand shape,
  location, direction (in progress) => exploit
  knowledge of „stroke"

# EMBRScript

```
BEGIN K_POSE
 TIME:1250
 HOLD:50

 BEGIN POSITION_CONSTRAINT
    BODY_GROUP:larm
    TARGET:0.3;-0.5;0.6
    JOINT:lhand
    OFFSET:0.0;0.0;0.0
 END
 BEGIN ORIENTATION_CONSTRAINT
    BODY_GROUP:larm
    NORMAL:Yaxis
    DIRECTION:0.0;-1.0;0.0
    JOINT:lhand
 END

END
```

- Pose: body configuration for a single time point (+ hold duration) defined by constraints like

  ➡ hand at a particular point in space

  ➡ hand shape, shoulder position

  ➡ facial expression, level of blushing

- Every constraint applies to part of the body

  ➡ Channels are inherent (arms, hands, face, shoulders, ...)

- Pose sequence: sequence of poses + start time

  ➡ a deliberate temporal segmentation

  ➡ design decision: we use sequences for glosses

```
BEGIN K_POSE_SEQUENCE
 CHARACTER:Alphonse
 START:390

 BEGIN K_POSE
 ...
 END

 BEGIN K_POSE
 ...
 END
END
```

# EMBRScript



Temporal „movement phase" markers allow synchronization & modification (e.g. drop preparation)

[Heloir, Kipp, 2009, 2010]

BehaviorBuilder tool to create and test EMBRScripts

# Sign Language Animation

- Attempt 1:
    - ➡ source video (human)
    - ➡ rotoscope (avatar)

# Why?

- Single sign disambiguation

  ➡ same manual movement, different meaning

  ➡ mouthing

  ➡ gaze, facial expression, pose narrow down possible meanings

- identify sentence topic

  ➡ interrogative facial expression / eyebrow raise

  ➡ pauses

  ➡ posture shift

- Face as fixation point

  ➡ allows parallel observance of face, mouth, hands, torso

  ➡ hard to do if face is static

# Sign Language Animation

- Attempt 1: *failure*

  ➡ source video (human)

  ➡ rotoscope (avatar)

- Attempt 2:

  ➡ source video (human)

  ➡ remake (human)

  ➡ rotoscope (avatar)

Gloss-wise transcription for utterance segmentation.

re-make

# BehaviorBuilder Extensions:
# Gloss-based Creation of Animation Sequences



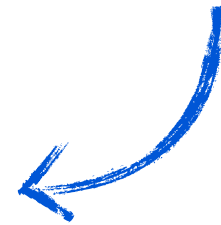Poses in current gloss

Glosses

EMBRScript

Current pose

original

re-make

avatar

# Lessons...

- SL is multimodal => change focus from manual gesture to ...
  - ➡ facial expression
  - ➡ mouthing
  - ➡ torso involvement
  - ➡ gaze

- Multimodality means
  - ➡ each modality as important as manual signs
  - ➡ explore synchronization

- Acceptability depends on
  - ➡ presence of style, personality, emotionality
  - ➡ prosody for information structure (topic) and segmentation
  - ➡ visual interest of the face => face as center of attention

- Good reliability test: Is it comprehensible?

# Conclusions

- We need motion capture!

- Sign language research needs you!

- Start looking at numbers instead of pixels...

# Thanks for listening!