THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

ISMAR 2016
September 19 - 23 | Merida, MEXICO

STEVENS
INSTITUTE of TECHNOLOGY
THE INNOVATION UNIVERSITY*
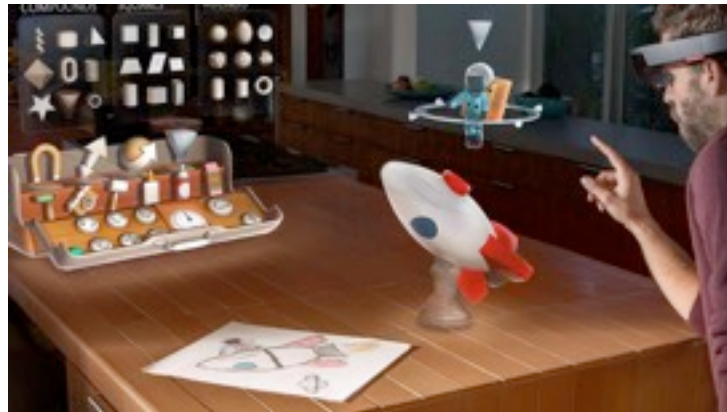
# Towards kHz 6-DoF Visual Tracking Using an Egocentric Cluster of Rolling Shutter Cameras

**Akash Bapat**[1], Enrique Dunn[1,2] & Jan-Michael Frahm[1],
UNC Chapel Hill, USA[1],
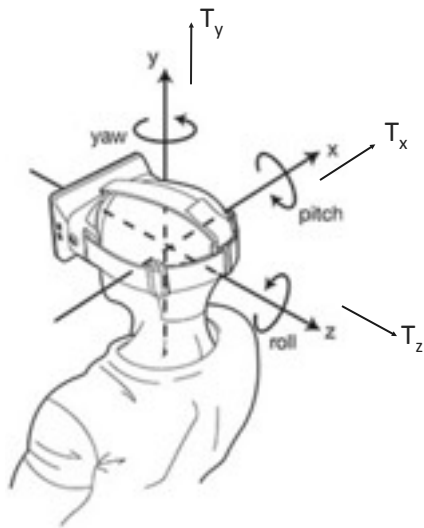Stevens Institute of Technology, USA[2]

# AR/VR System Components



Tracking ➤ Rendering ➤ Display

# AR/VR System Components



Find 6-DoF head pose[1]

---

Department of Computer Science    1.    LaValle etal, Head Tracking for Oculus Rift

Tuesday, September 20, 16

# Tracking Systems

Tuesday, September 20, 16

# Tracking Systems

Frequency

Complexity or cost

NDI

NDI Optotrak Certus[1]
- Tracking frequency ≈ 5kHz
- State of the art

Department of Computer Science
1. http://certus.ndigital.com
2. http://vrshoot.ru/sites/default/files/oculus-rift-dk2-cam.jpg
3. http://i.imgur.com/64LA9Kv.jpg
4. http://cdn.pocketnow.com/wp-content/uploads/2011/04/gyro.jpg
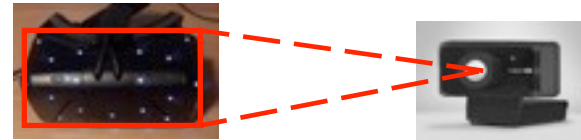5. LaValle etal, Head Tracking for Oculus Rift

Tuesday, September 20, 16

# Tracking Systems

Oculus Rift DK2



Orientation tracking
IMU at 1000 Hz [4,5]



Positional tracking [2,3]
Camera fps = 60 fps

Frequency

DK2

NDI

Complexity or cost

Department of Computer Science
1. http://certus.ndigital.com
2. http://vrshoot.ru/sites/default/files/oculus-rift-dk2-cam.jpg
3. http://i.imgur.com/64LA9Kv.jpg
4. http://cdn.pocketnow.com/wp-content/uploads/2011/04/gyro.jpg
5. LaValle etal, Head Tracking for Oculus Rift

Tuesday, September 20, 16

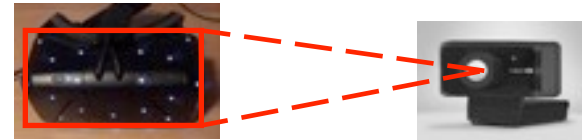# Tracking Systems



Oculus Rift DK2

Orientation tracking
IMU at 1000 Hz [4,5]

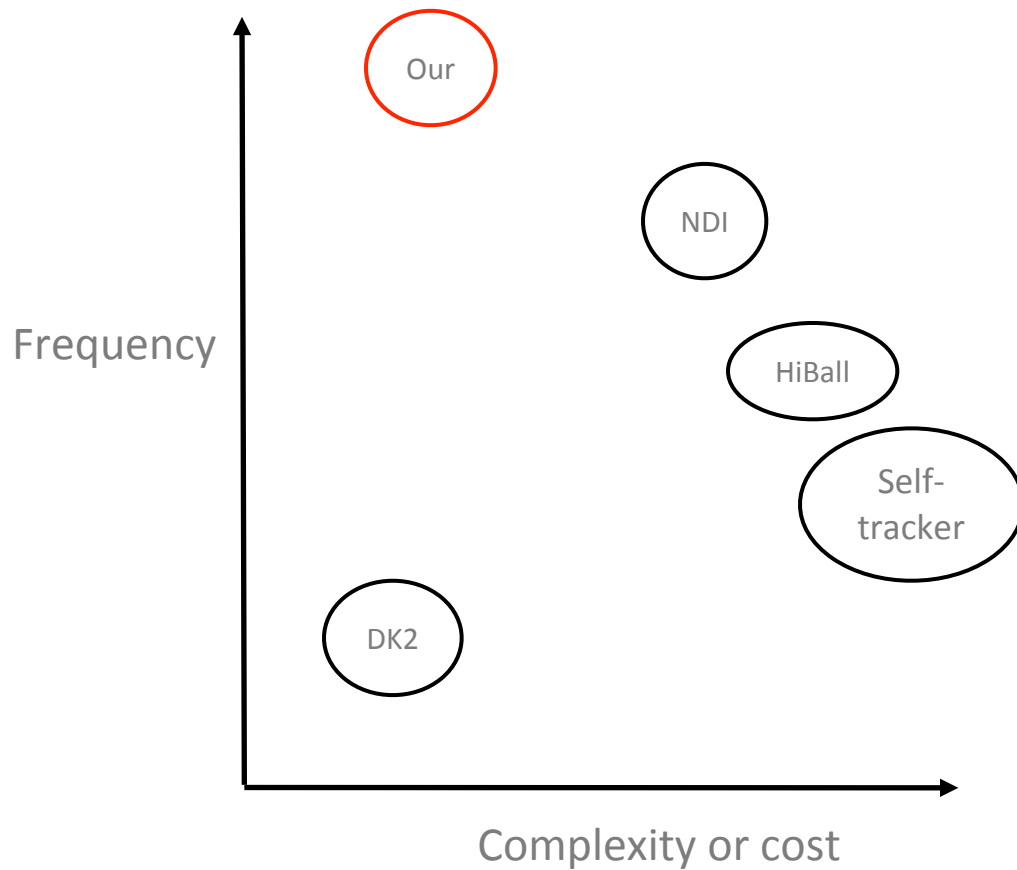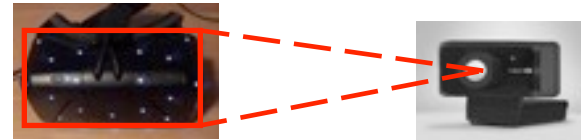Positional tracking [2,3]
Camera fps = 60 fps

# Tracking Systems



Frequency

Complexity or cost

Our

NDI

HiBall

Self-tracker

DK2

Oculus Rift DK2

Orientation tracking
IMU at 1000 Hz [4,5]

Positional tracking [2,3]
Camera fps = 60 fps

Department of Computer Science

1. http://certus.ndigital.com
2. http://vrshoot.ru/sites/default/files/oculus-rift-dk2-cam.jpg
3. http://i.imgur.com/64LA9Kv.jpg
4. http://cdn.pocketnow.com/wp-content/uploads/2011/04/gyro.jpg
5. LaValle etal, Head Tracking for Oculus Rift

Tuesday, September 20, 16

# Rolling Shutter



Rolling shutter capture[1]

Department of Computer Science

1. http://jasmcole.com/2014/10/12/rolling-shutters/

Tuesday, September 20, 16

# Rolling Shutter

- Row-by-row acquisition of linescan snapshots at slightly different times



Rolling shutter capture[1]

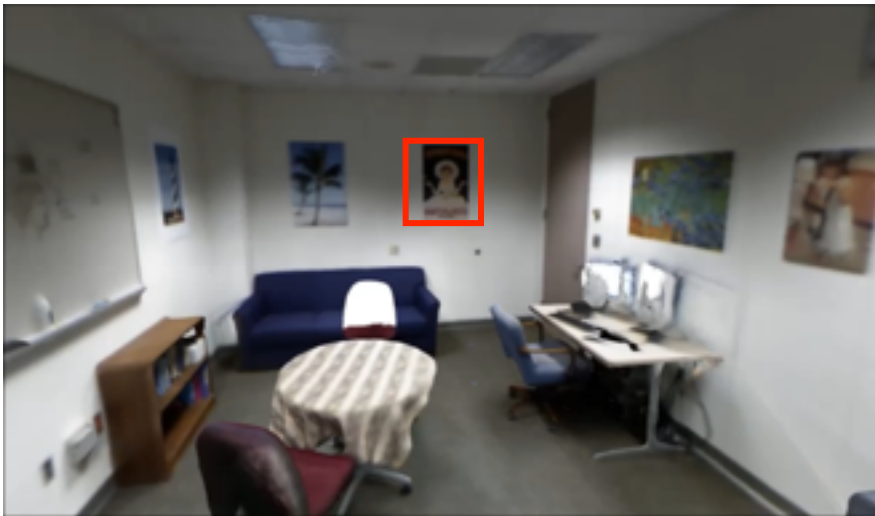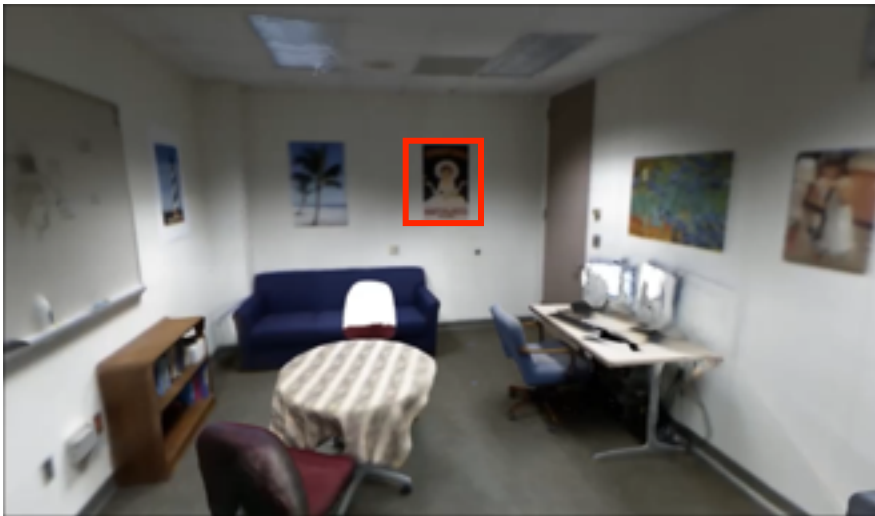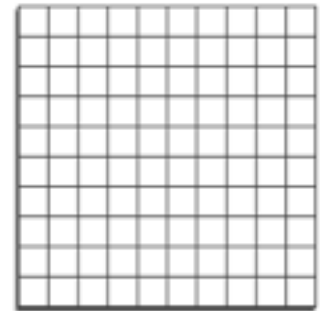1. http://jasmcole.com/2014/10/12/rolling-shutters/

Tuesday, September 20, 16

# Rolling Shutter

- Row-by-row acquisition of linescan snapshots at slightly different times

- Stream of row-images



Rolling shutter capture[1]

1. http://jasmcole.com/2014/10/12/rolling-shutters/

Tuesday, September 20, 16

# Rolling Shutter

- Row-by-row acquisition of linescan snapshots at slightly different times

- Stream of row-images



Rolling shutter capture[1]

1. http://jasmcole.com/2014/10/12/rolling-shutters/

Tuesday, September 20, 16

# Rolling Shutter

- Row-by-row acquisition of linescan snapshots at slightly different times

- Stream of row-images



Rolling shutter capture[1]

1. http://jasmcole.com/2014/10/12/rolling-shutters/

Tuesday, September 20, 16

# Stream of Rows

Rolling shutter[1]

- Frequency of row-samples
  - F = FPS * Height

$$= 120 * 720 > 80kHz$$

Tuesday, September 20, 16

# Stream of Rows

Rolling shutter[1]

- Frequency of row-samples
  - F = FPS * Height
  - =  120 * 720 > 80kHz

Tuesday, September 20, 16

# Our Tracker

Department of Computer Science

1. http://shop.gopro.com/hero4/hero4-black/CHDHX-401.html

Tuesday, September 20, 16

# Our Tracker

- Enabling component for rendering
  1. Zheng et al., 2014
  2. Lincoln et al., 2016

```
┌ ─ ─ ─ ─ ─ ─ ─ ─ ┐
  ┌──────────┐           ┌──────────┐
│ │ >30 kHz  │ │ ──────▶ │ 30 kHz   │
  │ tracker  │           │ Renderer │
│ └──────────┘ │         └──────────┘
└ ─ ─ ─ ─ ─ ─ ─ ─ ┘
```

Tuesday, September 20, 16

# Our Tracker

- Enabling component for rendering
  - 1. Zheng et al., 2014
  - 2. Lincoln et al., 2016

```
┌ ─ ─ ─ ─ ─ ─ ─ ─ ┐
|  ┌─────────┐    |      ┌─────────┐
|  │ >30 kHz │────┼─────▶│  30 kHz │
|  │ tracker │    |      │ Renderer│
|  └─────────┘    |      └─────────┘
└ ─ ─ ─ ─ ─ ─ ─ ─ ┘
```

GoPro[1]

Tuesday, September 20, 16

# Our Tracker

- Enabling component for rendering

    1. Zheng et al., 2014

    2. Lincoln et al., 2016



- Use commodity cameras



GoPro[1]

Department of Computer Science    1.    http://shop.gopro.com/hero4/hero4-black/CHDHX-401.html

Tuesday, September 20, 16

# Our Tracker

- Enabling component for rendering
    1. Zheng et al., 2014
    2. Lincoln et al., 2016



- Use commodity cameras



GoPro[1]

- Tracking frequency
    - Up to 80 kHz

- Break frame-rate barrier
    - Process each row of image

Department of Computer Science   1.   http://shop.gopro.com/hero4/hero4-black/CHDHX-401.html

Tuesday, September 20, 16

# Related Work



Artifact visualisation[1]

Artifact visualisation[1]

# Related Work

- ## Removing rolling shutter
    - Forssen et al., CVPR 2010
    - Track points using KLT tracker, estimate rotation
    - Parametrize intra-frame rotation as spline

Department of Computer Science   1.   http://jasmcole.com/2014/10/12/rolling-shutters/

Tuesday, September 20, 16

# Related Work

- ## Removing rolling shutter

  - Forssen et al., CVPR 2010

  - Track

  - Para

Tuesday, September 20, 16

# Related Work



Artifact visualisation[1]

- ## Removing rolling shutter
  - Forssen et al., CVPR 2010
  - Track points using KLT tracker, estimate rotation
  - Parametrize intra-frame rotation as spline



- ## Geometric problems
  - Multi-view stereo: Saurer et al., ICCV 2013
  - Adapt plane sweep stereo for rolling shutter

Department of Computer Science   1.   http://jasmcole.com/2014/10/12/rolling-shutters/

Tuesday, September 20, 16

# Related Work



Artifact visualisation[1]

- ## Removing rolling shutter
  - Forssen et al., CVPR 2010
  - Track points using KLT tracker, estimate rotation
  - Parametrize intra-frame rotation as spline



- ## Geometric problems
  - Multi-view stereo: Saurer et al., ICCV 2013
  - Adapt plane sweep stereo for rolling shutter



GT depth    GS error    RS error

- ## Velocity estimation
  - Ait-Aider et al., ICVS 2006
  - Solve for pose and velocity using bundle adjustment
  - Use 2D-3D correspondence, similar to camera calibration

Tuesday, September 20, 16

# Related Work

- Removing rolling shutter



- Solve for pose and velocity using bundle adjustment
- Use 2D-3D correspondence, similar to camera calibration

Department of Computer Science   1.   http://jasmcole.com/2014/10/12/rolling-shutters/

Tuesday, September 20, 16

# Related Work



Artifact visualisation[1]

- ## Removing rolling shutter
  - Forssen et al., CVPR 2010
  - Track points using KLT tracker, estimate rotation
  - Parametrize intra-frame rotation as spline



- ## Geometric problems
  - Multi-view stereo: Saurer et al., ICCV 2013
  - Adapt plane sweep stereo for rolling shutter



GT depth     GS error     RS error

- ## Velocity estimation
  - Ait-Aider et al., ICVS 2006
  - Solve for pose and velocity using bundle adjustment
  - Use 2D-3D correspondence, similar to camera calibration

Tuesday, September 20, 16

# Approach Overview

Tuesday, September 20, 16

# Approach Overview

Camera Cluster

Tuesday, September 20, 16

# Approach Overview

Camera Cluster

Spatial Stereo

Tuesday, September 20, 16

# Approach Overview

Department of Computer Science

Tuesday, September 20, 16

# Approach Overview

Camera Cluster

Spatial Stereo

Temporal Stereo

Department of Computer Science

Tuesday, September 20, 16

# Approach Overview

Estimation Framework

Camera Cluster

↓

Spatial Stereo

↓

Temporal Stereo

Tuesday, September 20, 16

# Approach Overview



Our approach

Department of Computer Science

Tuesday, September 20, 16

# Approach Overview



Cluster

Estimation Framework

Camera Cluster

Spatial Stereo

Temporal Stereo

Stream of rows

Input

Our approach

Department of Computer Science

Tuesday, September 20, 16

# Approach Overview



Cluster

Input

Stream of rows

Estimation Framework

Camera Cluster

Spatial Stereo

Temporal Stereo

Our approach

High frequency 6DoF tracking

Output

Department of Computer Science

Tuesday, September 20, 16

# Our Cluster

GoPro[1]

1. http://shop.gopro.com/hero4/hero4-black/CHDHX-401.html
2. http://users.erols.com/njastro/faas/image001/CCD2.jpg

Tuesday, September 20, 16

# Our Cluster

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

GoPro[1]

Dynamic sensor with
periodic movement[2]

Tuesday, September 20, 16

# Our Cluster

GoPro[1]



Dynamic sensor with periodic movement[2]

Department of Computer Science

1. http://shop.gopro.com/hero4/hero4-black/CHDHX-401.html
2. http://users.erols.com/njastro/faas/image001/CCD2.jpg

Tuesday, September 20, 16

# Our Cluster

GoPro[1]

Dynamic sensor with
periodic movement[2]

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

Department of Computer Science

1.  http://shop.gopro.com/hero4/hero4-black/CHDHX-401.html
2.  http://users.erols.com/njastro/faas/image001/CCD2.jpg

Tuesday, September 20, 16

# Our Cluster

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

GoPro[1]

Dynamic sensor with
periodic movement[2]

Small vertical FoV



10

Department of Computer Science

1.  http://shop.gopro.com/hero4/hero4-black/CHDHX-401.html
2.  http://users.erols.com/njastro/faas/image001/CCD2.jpg

Tuesday, September 20, 16

# Our Cluster

# Our Cluster

🔴 N stereo-pairs of rolling shutter camera, N=5

- Precalibrated intrinsic and extrinsics
- Temporal sync
- Known history of motion

🟢 Hi-Ball for ground truth

# Our Cluster

🔴 N stereo-pairs of rolling shutter camera, N=5

- Precalibrated intrinsic and extrinsics
- Temporal sync
- Known history of motion

Ⓗ Hi-Ball for ground truth

Geometry of cluster

Department of Computer Science

Tuesday, September 20, 16

# Motion Model

3D point

$$X^{t+1}_{cam_i} = M_{cam_i} V^{-1} \boxed{X^t_{cam_i}}$$

3D point X

Tuesday, September 20, 16

# Motion Model

3D point

$$X^{t+1}_{cam_i} = M_{cam_i} \boxed{X^{t}_{cam_i}}$$

3D point X

Tuesday, September 20, 16

# Motion Model

3D point

$$X^{t+1}_{cam_i} = M_{cam_i} \quad X^t_{cam_i}$$

Motion of camera

3D point X

Department of Computer Science

Tuesday, September 20, 16

# Motion Model

Moved 3D point

3D point

$$X_{cam_i}^{t+1} = M_{cam_i} \; X_{cam_i}^{t}$$

Motion of camera

3D point X

Department of Computer Science

Tuesday, September 20, 16

# Motion Model

Moved 3D point

3D point

$$X_{cam_i}^{t+1} = M_{cam_i} \quad X_{cam_i}^{t}$$

$$X_{cam_i}^{t+1} = V_i \, M_{cluster} \, V_i^{-1} \, X_{cam_i}^{t}$$

$v_i$

$cam_i$

RS cluster

Tuesday, September 20, 16

# Linearized Motion Model

$$X^{t+1}_{cam_i} = M_{cam_i} \qquad X^t_{cam_i}$$

$$X^{t+1}_{cam_i} = V_i \, M_{cluster} \, V_i^{-1} \, X^t_{cam_i}$$

$$M_{cluster} = \begin{bmatrix} R(\theta_x, \theta_y, \theta_z) & T \\ 0 & 1 \end{bmatrix}$$

Department of Computer Science

15

# Linearized Motion Model

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

$$X^{t+1}_{cam_i} = M_{cam_i} X^t_{cam_i}$$

$$X^{t+1}_{cam_i} = V_i \, M_{cluster} \, V_i^{-1} \, X^t_{cam_i}$$

$$M_{cluster} = \begin{bmatrix} R(\theta_x, \theta_y, \theta_z) & T \\ 0 & 1 \end{bmatrix} \implies dM = \begin{bmatrix} 1 & -\theta_z & \theta_y & \delta T_x \\ \theta_z & 1 & -\theta_x & \delta T_y \\ -\theta_y & \theta_x & 1 & \delta T_z \\ 0 & 0 & 0 & 1 \end{bmatrix}$$
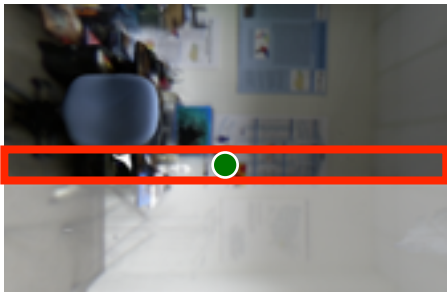
Tuesday, September 20, 16
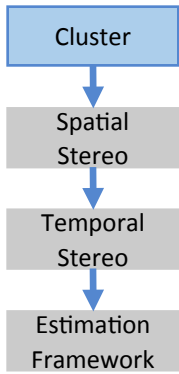
# Linearized Motion Model

$$X_{cam_i}^{t+1} \quad = \quad M_{cam_i} \quad X_{cam_i}^t$$

$$X_{cam_i}^{t+1} \quad = V_i \, M_{cluster} \, V_i^{-1} \, X_{cam_i}^t$$

$$M_{cluster} = \begin{bmatrix} R(\theta_x, \theta_y, \theta_z) & T \\ 0 & 1 \end{bmatrix} \implies dM = \begin{bmatrix} 1 & -\theta_z & \theta_y & \delta T_x \\ \theta_z & 1 & -\theta_x & \delta T_y \\ -\theta_y & \theta_x & 1 & \delta T_z \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Department of Computer Science

15

Tuesday, September 20, 16

# Linearized Motion Model

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

$$X^{t+1}_{cam_i} = M_{cam_i} X^t_{cam_i}$$

$$X^{t+1}_{cam_i} = V_i M_{cluster} V_i^{-1} X^t_{cam_i}$$

$$M_{cluster} = \begin{bmatrix} R(\theta_x, \theta_y, \theta_z) & T \\ 0 & 1 \end{bmatrix} \implies dM = \begin{bmatrix} 1 & -\theta_z & \theta_y & \delta T_x \\ \theta_z & 1 & -\theta_x & \delta T_y \\ -\theta_y & \theta_x & 1 & \delta T_z \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Department of Computer Science

15

Tuesday, September 20, 16

# Motion Estimation

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

$$X_{cam_i}^{t+1} = V_i \; M_{cluster} \; V_i^{-1} \; \boxed{X_{cam_i}^t}$$

$$X_{cam_i}^t = [0 \quad y \quad d \quad 1]^T$$

Department of Computer Science

Tuesday, September 20, 16

# Motion Estimation

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

$$X^{t+1}_{cam_i} \;=\; V_i \; M_{cluster} \; V_i^{-1} \; \boxed{X^t_{cam_i}}$$

$$X^t_{cam_i} = [0 \quad y \quad d \quad 1]^T$$

Stereo in space

Tuesday, September 20, 16

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

# Motion Estimation

$$X_{cam_i}^{t+1} = [\Delta x \quad y + \Delta y \quad d + \Delta d \quad 1]^T$$

$$\boxed{X_{cam_i}^{t+1}} = V_i \ M_{cluster} \ V_i^{-1} \ X_{cam_i}^t$$

Tuesday, September 20, 16

# Motion Estimation

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

$$X_{cam_i}^{t+1} = [\Delta x \quad y + \Delta y \quad d + \Delta d \quad 1]^T$$

$$\boxed{X_{cam_i}^{t+1}} = V_i \ M_{cluster} \ V_i^{-1} \ X_{cam_i}^t$$

Tuesday, September 20, 16

# Motion Estimation

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

$$X_{cam_i}^{t+1} = [\Delta x \quad y + \Delta y \quad d + \Delta d \quad 1]^T$$

$$\boxed{X_{cam_i}^{t+1}} = V_i \; M_{cluster} \quad V_i^{-1} \quad X_{cam_i}^{t}$$



---

Department of Computer Science

Tuesday, September 20, 16

# Stereo Estimation

Department of Computer Science

Tuesday, September 20, 16

# Stereo Estimation

- ## Spatial stereo $s_d$
  - Use stereo-pair of cameras
  - Measure pixel disparity

- ## Temporal stereo $s_t$
  - Measure small shifts in pixel
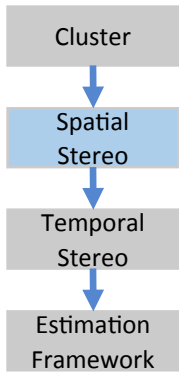  - Compare row at time t and t-k

t1 $\longrightarrow$ t2

Department of Computer Science

Tuesday, September 20, 16

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

# Stereo Estimation

- # Spatial stereo $s_d$

  - ## Use stereo-pair of cameras
  - ## Measure pixel disparity

- # Temporal stereo $s_t$

  - ## Measure small shifts in pixel
  - ## Compare row at time t and t-k

t1 $\longrightarrow$ t2

Tuesday, September 20, 16

# Binary Row Descriptor

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

# Binary Row Descriptor

Visualization of row-image
(296th)

# Binary Row Descriptor

Visualization of row-image (296th)

$$\frac{d^2 G(x, 0, \sigma)}{dx^2} * row$$
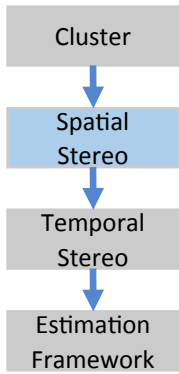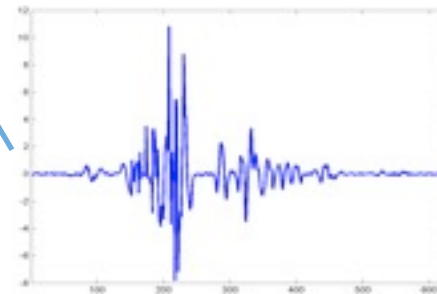
2nd derivative

# Binary Row Descriptor

Visualization of row-image (296th)

Sign

$$\frac{d^2 G(x, 0, \sigma)}{dx^2} * row$$

2nd derivative

Tuesday, September 20, 16

Cluster

Spatial
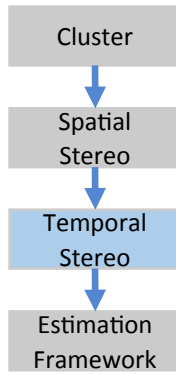Stereo

Temporal
Stereo

Estimation
Framework

# Spatial Stereo : Measure Disparity

- Compare binary descriptor of rows of stereo cameras

- Fast hamming cost matching

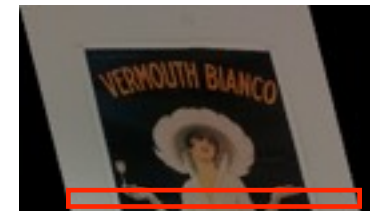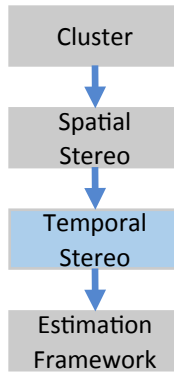- $depth = \dfrac{focal * baseline}{disparity}$

Stereo in space

P (x,y,d)

$x_L$

$x_R$

focal length

baseline

Stereo Pair

Row samples of stereo pair

1 0 1 0 1 1 1 1 0 0       0 1 0 1 0 1 1 1 1 0

Binary descriptor

$S_d$

Current frames

# Temporal Stereo: Measure Shift $s_t$

- Different regions in space are captured at different timestamps

- Stereo in time
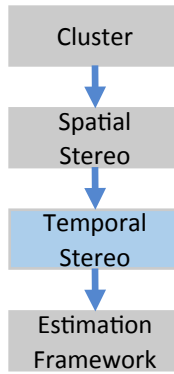  - Need snaps of same space but at different timestamps

Tuesday, September 20, 16

# Temporal Stereo: Measure Shift $s_t$

- Different regions in space are captured at different timestamps

- Stereo in time
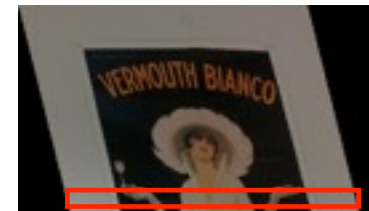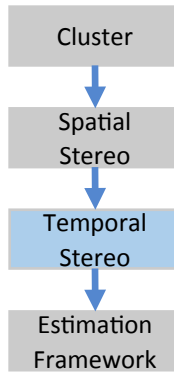    - Need snaps of same space but at different timestamps



Current time

---

Tuesday, September 20, 16

# Temporal Stereo: Measure Shift $s_t$

- Different regions in space are captured at different timestamps

- Stereo in time
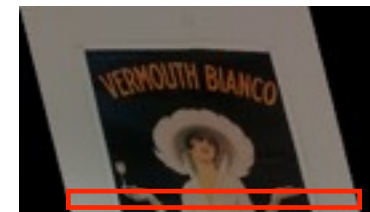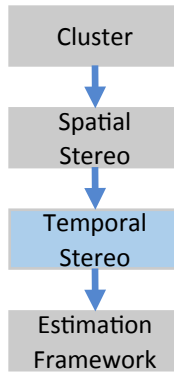  - Need snaps of same space but at different timestamps



Previous Frame

Current time

Department of Computer Science
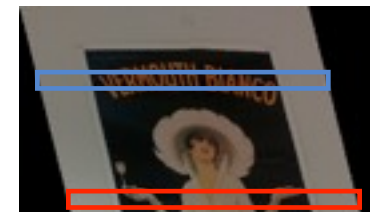
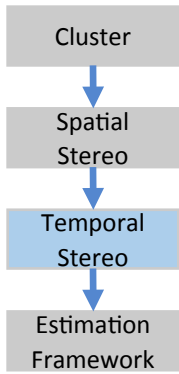Tuesday, September 20, 16

# Temporal Stereo: Measure Shift $s_t$

- Different regions in space are captured at different timestamps

- Stereo in time
  - Need snaps of same space but at different timestamps



Previous Frame

Known [R |t] per row

Current time

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

# Temporal Stereo: Measure Shift $s_t$

- Different regions in space are captured at different timestamps

- Stereo in time
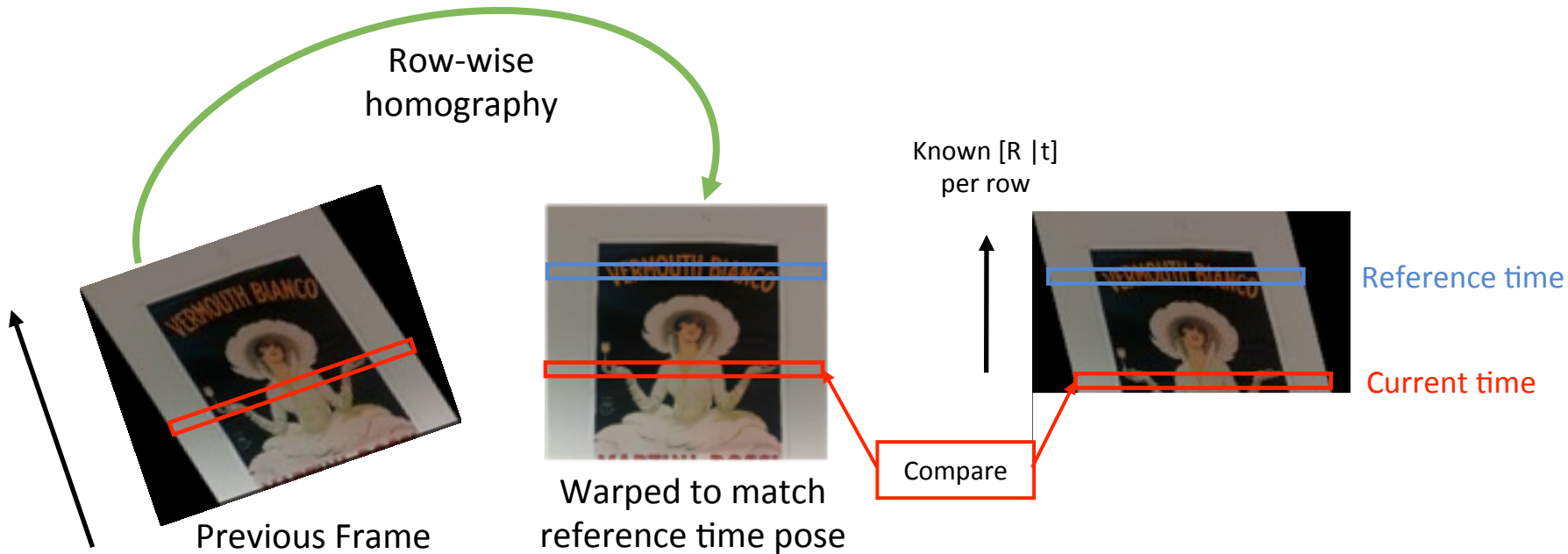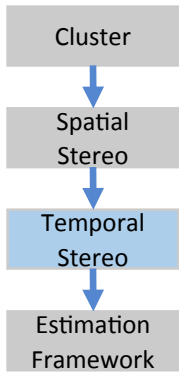    - Need snaps of same space but at different timestamps



Previous Frame

Known [R |t]
per row

Reference time

Current time

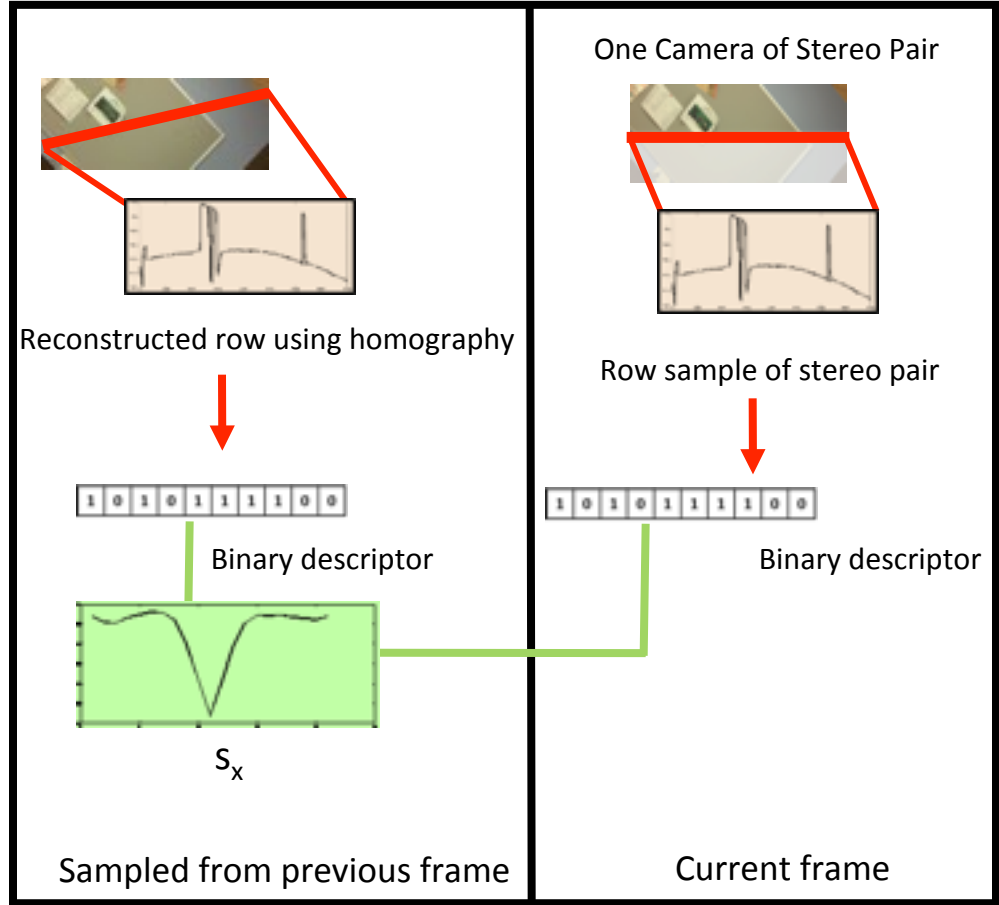Tuesday, September 20, 16

# Temporal Stereo: Measure Shift $s_t$

- Different regions in space are captured at different timestamps

- Stereo in time
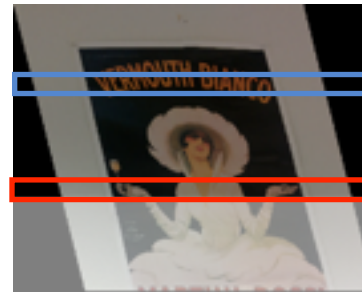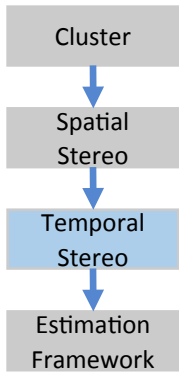  - Need snaps of same space but at different timestamps



Row-wise homography

Known [R |t] per row

Previous Frame

Warped to match reference time pose

Compare

Reference time

Current time

# Temporal Stereo: Measure Shift $s_t$

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

- Reconstruct row using per- row homography

- Use binary descriptors and hamming distance



One Camera of Stereo Pair

Reconstructed row using homography

Row sample of stereo pair

1 0 1 0 1 1 1 1 0 0

1 0 1 0 1 1 1 1 0 0

Binary descriptor

Binary descriptor

$s_x$

Sampled from previous frame

Current frame

Tuesday, September 20, 16

# Adaptive Reference

Reference row: t-k

Current Row (t)

Current frame

Estimation reliability

0

k

23

# Adaptive Reference

- ## Leave sufficient motion
  - Row-to-row motion
  - Interpolation & pixel measurement noise
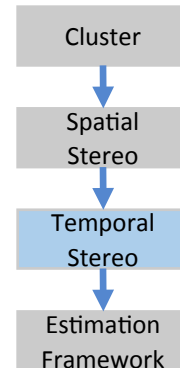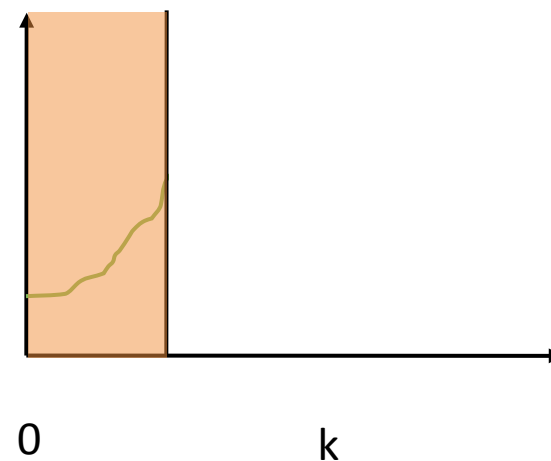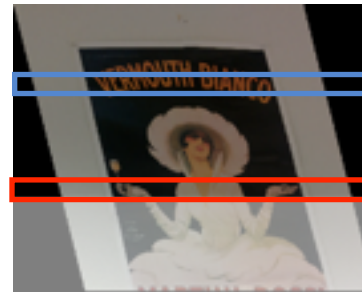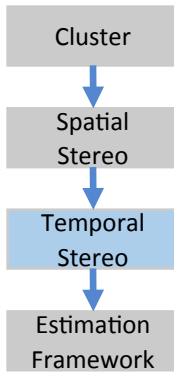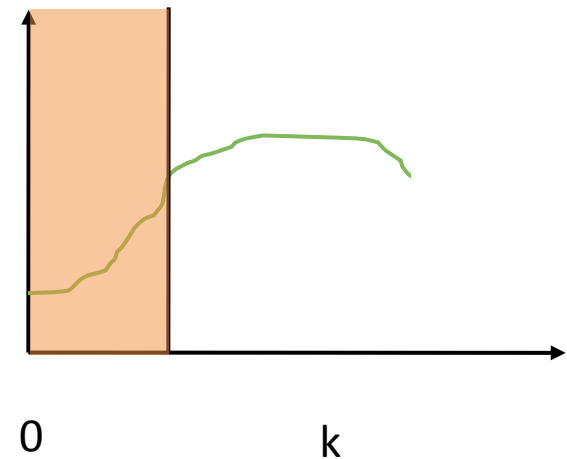


Reference row: t-k

Current Row (t)

Current frame

- ## Satisfy small-motion assumption
  - Reference should not be far away
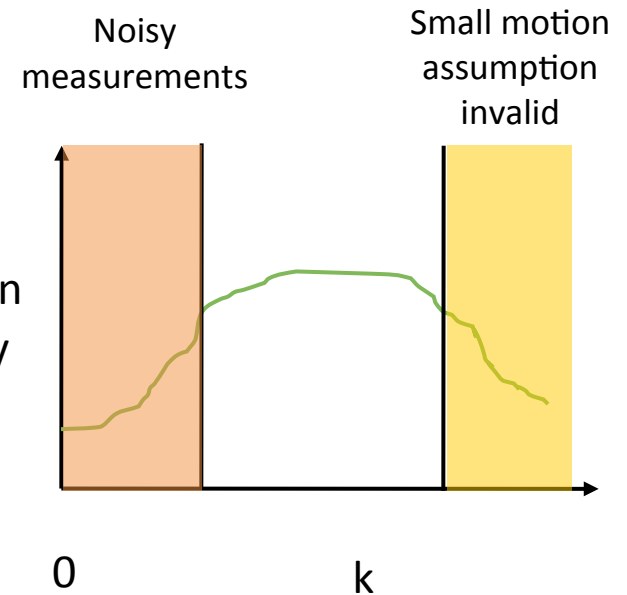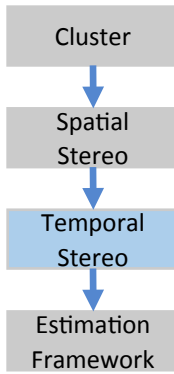
Noisy measurements



Estimation reliability

0          k

Department of Computer Science

Tuesday, September 20, 16

# Adaptive Reference

- Leave sufficient motion
  - Row-to-row motion
  - Interpolation & pixel measurement noise

Reference row: t-k

Current Row (t)

Current frame

Noisy measurements

- Satisfy small-motion assumption
  - Reference should not be far away

Estimation reliability

0                    k

Department of Computer Science

23

Tuesday, September 20, 16

# Adaptive Reference

- Leave sufficient motion
  - Row-to-row motion
  - Interpolation & pixel measurement noise

Reference row: t-k

Current Row (t)

Current frame

- Satisfy small-motion assumption
  - Reference should not be far away

Noisy measurements

Small motion assumption invalid

Estimation reliability

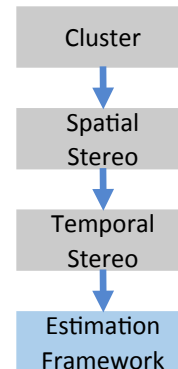0　　　　　k

23

# Confidence Scores

High confidence | Low Confidence

1. Quality of minimum: $c_{PKR}$
   - Is the valley unique?
   - Use Peak ratio (PKR)[1]

2. Temporal consistency : $c_t$
   - Consistent shifts in time
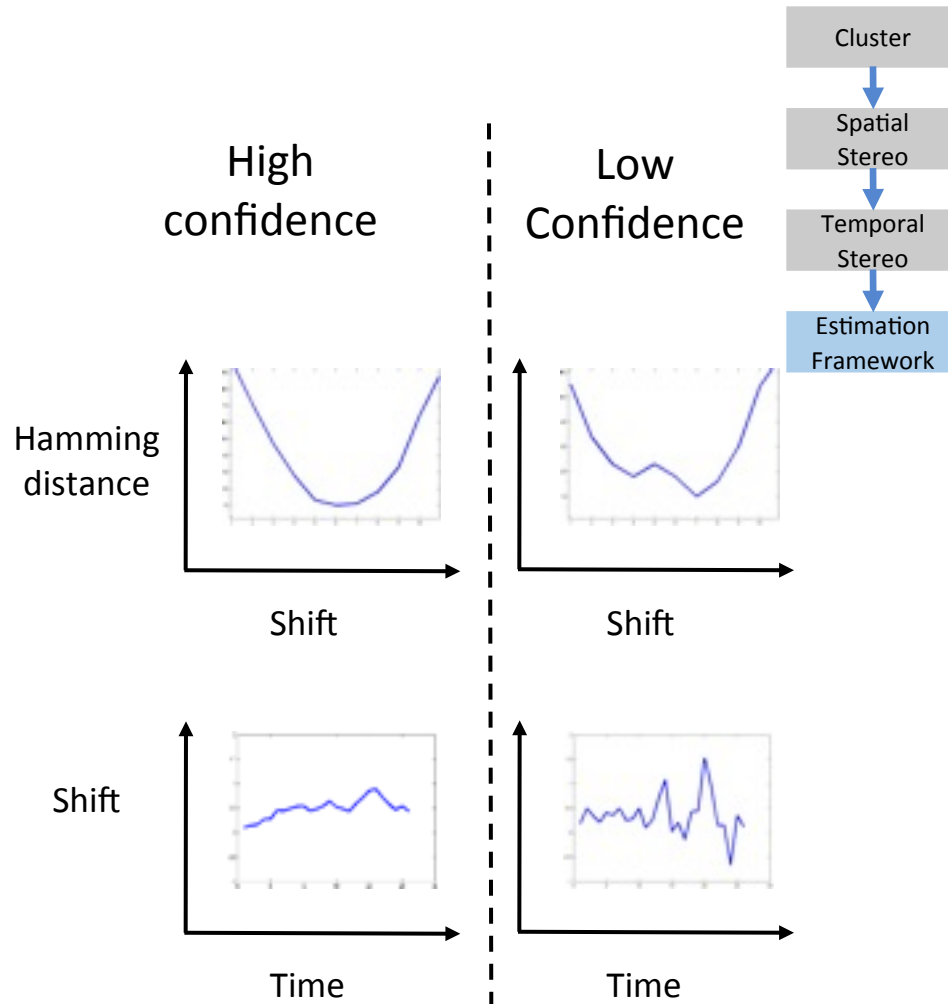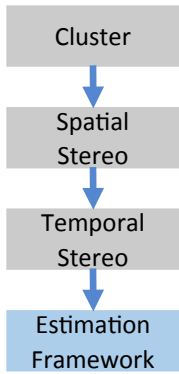   - Penalize sudden changes

3. $c_{i,t} = c_{PKR} * c_t$

Department of Computer Science

1. X. Hu and P. Mordohai. A quantitative evaluation of confidence measures for stereo vision, PAMI 2012

Tuesday, September 20, 16

# Confidence Scores

**High confidence** | **Low Confidence**

1. Quality of minimum: $c_{PKR}$
   - Is the valley unique?
   - Use Peak ratio (PKR)[1]

Hamming distance



Shift     Shift

2. Temporal consistency : $c_t$
   - Consistent shifts in time
   - Penalize sudden changes

Shift



Time     Time

3. $c_{i,t} = c_{PKR} * c_t$

$$\mathbf{C}(t) = \begin{bmatrix} c_{1,t} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & c_{n,t} \end{bmatrix}$$

24

1. X. Hu and P. Mordohai. A quantitative evaluation of confidence measures for stereo vision, PAMI 2012

Tuesday, September 20, 16

# Motion Estimation

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

$$X_{cam_i}^{t+1} = [\Delta x \quad y + \Delta y \quad d + \Delta d \quad 1]^T$$

$$X_{cam_i}^{t+1} = V_i \ M_{cluster} \ V_i^{-1} \ X_{cam_i}^t$$

Stereo in time

$$X_{cam_i}^t = [0 \quad y \quad d \quad 1]^T$$

Stereo in space

Tuesday, September 20, 16

Cluster

Spatial
Stereo

Temporal
Stereo

Estimation
Framework

# Weighted Linear System

- $X^{t+1}_{cam_i} = V_i\, M_{cluster}\, V_i^{-1} X^t_{cam_i}$

  - One equation from each camera
  - Use more cameras for robustness

$$C\,A \begin{bmatrix} \delta T_x \\ \delta T_y \\ \delta T_z \\ \theta_x \\ \theta_y \\ \theta_z \end{bmatrix} = C\,B$$

Department of Computer Science

26

Tuesday, September 20, 16

# Approach Summary

Department of Computer Science

Tuesday, September 20, 16

# Approach Summary



Stream of rows

Input

Tuesday, September 20, 16

# Approach Summary



Our approach

Input

Tuesday, September 20, 16

# Approach Summary

Tuesday, September 20, 16

# Simulator

- Developed in OpenGL+Qt and Unity3D
  - Support for Hi-Ball tracker data
  - Motion Blur

Tuesday, September 20, 16

# Simulator

- Developed in OpenGL+Qt and Unity3D
  - Support for Hi-Ball tracker data
  - Motion Blur

Department of Computer Science

Tuesday, September 20, 16

# Experiments

- Camera specs:
  - 640 × 480 pixels
  - FoV = 60°
  - Stereo pairs = 10
  - 120 fps



- Quantify errors in terms of display pixel errors
  - Display specs : HTC Vive
    - 120° FoV
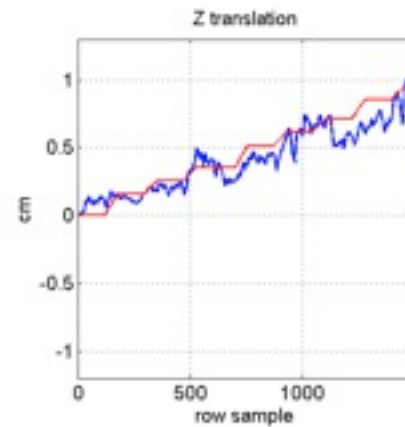    - 1080 × 1200 pixels  per eye
    - Point 1m in front

1.Mingsong Duo et al, Exploring high-level plane primitives for indoor 3d reconstruction with a hand-held rgb-d camera, ACCV 2012

# Display Pixel Error



Pixels shown in HMD

Tuesday, September 20, 16

# Display Pixel Error



Pixels shown in HMD

Department of Computer Science

Tuesday, September 20, 16
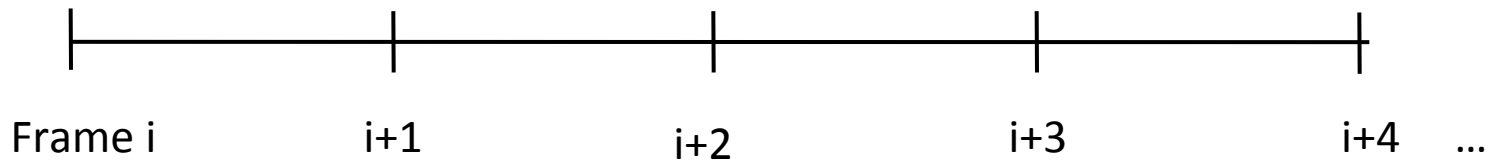
# Experiments : Real motion[1] in Synthetic Room



f= 56.7kHz

Department of Computer Science    1. Hi-Ball tracker data

Tuesday, September 20, 16

# Experiments : Real motion[1] in Synthetic Room

f= 56.7kHz



Frame i      i+1      i+2      i+3      i+4   ...

Tuesday, September 20, 16

# Experiments: Synthetic Large Motion



f= 56.7kHz
v= 1.4 m/s
ω= 120 deg/s

Tuesday, September 20, 16

# Experiments: Synthetic Large Motion

f= 56.7kHz

Department of Computer Science

Tuesday, September 20, 16

# Experiments : Synthetic  Extreme Motion



f= 56.7kHz

v= 1.4 m/s

ω= 500 deg/s

Tuesday, September 20, 16

# Experiments : Synthetic Extreme Motion

f= 56.7kHz
v= 1.4 m/s
ω= 500 deg/s



Sampling not possible

Tuesday, September 20, 16

# Results for Real Imagery



f= 80.4kHz

Department of Computer Science

Tuesday, September 20, 16

# Conclusion

- High frequency visual tracker
    - Up to 80 kHz
    - Off-the-shelf cameras

High Frequency

Tuesday, September 20, 16

# Conclusion

- High frequency visual tracker
  - Up to 80 kHz
  - Off-the-shelf cameras

Department of Computer Science

Tuesday, September 20, 16

# Conclusion

- High frequency visual tracker
    - Up to 80 kHz
    - Off-the-shelf cameras

Department of Computer Science

Tuesday, September 20, 16

# Conclusion

- High frequency visual tracker
  - Up to 80 kHz
  - Off-the-shelf cameras

Tuesday, September 20, 16

# Conclusion

- High frequency visual tracker
  - Up to 80 kHz
  - Off-the-shelf cameras

Tuesday, September 20, 16

# Thank you!

Contact: akash@cs.unc.edu

# Camera tracking : The spectrum

| Global Shutter | Rolling Shutter as 1-D sensor | 1-D sensor |
|---|---|---|



1-D camera[1]

| Global Shutter | Rolling Shutter as 1-D sensor | 1-D sensor |
|---|---|---|
| + No distortion | − Distortion | + No distortion |
| + Drift correction | + Drift correction | − No drift correction |
| − High cost | + Cheap | + High Cost |
| − Low fps ≈ 200Hz | + High fps, 56 kHz | + High fps, till 87 kHz |
| − Higher noise | + Low noise | − Higher noise |

Tuesday, September 20, 16

Department of Computer Science

Tuesday, September 20, 16