Shape and Motion under Varying Illumination: Unifying Structure from Motion, Photometric Stereo, and Multi-view Stereo

Li Zhang* **Brian Curless*** *University of Washington, Seattle, WA, USA







Introduction

Goal

Dense shape reconstruction of moving objects under varying illumination from a single video.

Example A hand-held figurine rotating in front of a fixed camera under static lighting.

Standard methods

- Structure from Motion – only feature points
- ♦ Multi-view Stereo - only constant lighting
- Photometric Stereo -
- only static objects

Our solution

Using both spatial and temporal brightness variations

spatial brightness variation temporal brightness variation photometric cue motion cue surface position surface orientation

dense surface reconstruction in both textured and textureless regions

Contributions

- In the dense structure from motion
- stereo matching under lighting changes

Aaron Hertzmann Steven M. Seitz*** **University of Toronto, Toronto, ON, Canada





Mathematical Formulation **Optical flow under varying illumination**

Assume a Lambertian object moving rigidly in front of an orthographic camera under static distant light.



Let \mathbf{l}_{t} and \mathbf{b}_{t} be the directional and ambient light at frame t, and $\mathbf{n}_{\rm p}$ and $\alpha_{\rm p}$ be the normal and albedo of the point p, then

Image Formation:

Irradiance ratio:

$$I_{t}(\mathbf{x}_{t,p}) = \boldsymbol{\alpha}_{p} \cdot (\mathbf{l}_{t}^{T} \mathbf{n}_{p} + \mathbf{b}_{t})$$
$$\frac{I_{t}(\mathbf{x}_{t,p})}{I_{0}(\mathbf{x}_{0,p})} = \frac{\mathbf{l}_{t}^{T} \mathbf{n}_{p} + \mathbf{b}_{t}}{\mathbf{l}_{0}^{T} \mathbf{n}_{p} + \mathbf{b}_{0}} = \gamma_{t,p}$$

Brightness-varying flow: $I_t(\mathbf{x}_{t,p}) = \gamma_{t,p} \cdot I_0(\mathbf{x}_{0,p})$

Multi-point multi-frame optical flow can be computed by minimizing: $\Phi(\{\mathbf{x}_{t,p}, \gamma_{t,p}\}) = \sum (I_t(\mathbf{x}_{t,p}) - \gamma_{t,p} \cdot I_0(\mathbf{x}_{0,p}))^2$

Constraints on brightness-varying flow

Rank 3 constraint on {x_{t,p}, y_{t,p}} [Tomasi&Kanade90]:

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_{1,1} \cdots \mathbf{X}_{1,P} \\ \vdots & \vdots & \vdots \\ \mathbf{X}_{T,1} \cdots & \mathbf{X}_{T,P} \end{bmatrix} = \begin{bmatrix} \vdots \\ \mathbf{r}_{xt}^{T} \\ \vdots \end{bmatrix} \begin{bmatrix} \cdots \mathbf{s}_{p} \cdots \end{bmatrix} + \begin{bmatrix} \mathbf{0}_{x1} \cdots \mathbf{0}_{x1} \\ \vdots & \vdots \\ \mathbf{0}_{xT} \cdots \mathbf{0}_{xT} \end{bmatrix} = \mathbf{R}_{x}\mathbf{S} + \mathbf{O}_{x}$$

 $\mathbf{Y} = \mathbf{R}_{v}\mathbf{S} + \mathbf{O}_{v}$, similarly.

Rank 4 constraint on $\{\gamma_{t,p}\}$ [Basri&Jacob01]:

$$\Gamma = \begin{bmatrix} \gamma_{1,1} \cdots \gamma_{1,P} \\ \vdots & \vdots & \vdots \\ \gamma_{T,1} \cdots & \gamma_{T,P} \end{bmatrix} = \begin{bmatrix} \vdots & \vdots \\ \mathbf{l}_t^T & \mathbf{b}_t \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} \cdots & \frac{\mathbf{n}_p}{\beta_p} & \cdots \\ \cdots & \frac{1}{\beta_p} & \cdots \\ \cdots & \frac{1}{\beta_p} & \cdots \end{bmatrix} = \mathbf{L} \widetilde{\mathbf{N}}$$

where $\beta_p = \mathbf{l}_0^T \mathbf{n}_p + \mathbf{b}_0$.

Consistency constraint on **S** and **N**: $\frac{\delta S}{\delta x} \perp N$, $\frac{\delta S}{\delta y} \perp N$.

Multi-point multi-frame brightness-varying optical flow min $\Phi(\mathbf{X}, \mathbf{Y}, \Gamma)$ 22 20

s.t.
$$\mathbf{X} = \mathbf{R}_{\mathbf{X}}\mathbf{S} + \mathbf{O}_{\mathbf{X}}, \ \mathbf{Y} = \mathbf{R}_{\mathbf{y}}\mathbf{S} + \mathbf{O}_{\mathbf{y}}, \ \Gamma = \mathbf{L}\mathbf{N}, \ \frac{\delta\mathbf{S}}{\delta\mathbf{x}} \perp \mathbf{N}, \ \frac{\delta\mathbf{S}}{\delta\mathbf{y}} \perp \mathbf{N}.$$

Reconstruction of the figurine



Reconstruction using only motion cue

The coarse-to-fine reconstruction, using both motion and photometric cues

Reconstruction of a nearly textureless object



however, the temporal brightness variations uniquely determine surface orientation.

Conclusion

- Motion cue constrains 3D positions inaccurately for low textured pixels
- Photometric cue reveals normals accurately even for moving scenes, despite the noisy motion estimation
- Combing both cues recovers moving shape densely

Our formulation extends [Irani99] and applies to features, edges, and textureless regions. Our formulation also subsumes Structure from Motion, Multi-view Stereo, and Photometric Stereo as special cases:

	Known	Unnown
Structure from Motion	X , Y	$\mathbf{R}_{\mathrm{X}}, \mathbf{O}_{\mathrm{X}}, \mathbf{R}_{\mathrm{y}}, \mathbf{O}_{\mathrm{y}}, \mathbf{S}$
Multi-view Stereo	$\mathbf{R}_{\mathrm{X}}, \mathbf{O}_{\mathrm{X}}, \mathbf{R}_{\mathrm{y}}, \mathbf{O}_{\mathrm{y}}, \Gamma=1$	S
Photometric Stereo	Γ, constant X, Y	N, L

Reconstruction algorithm

Initialization Track features, compute \mathbf{R}_{x} , \mathbf{O}_{x} , \mathbf{R}_{y} , \mathbf{O}_{y} , estimate feature normals, and initialize L.

Step1 Fix \mathbf{R}_x , \mathbf{O}_x , \mathbf{R}_v , \mathbf{O}_v , and \mathbf{L} , and update \mathbf{S} and \mathbf{N} ; Step2 Correct S with large uncertainties by integrating N; Step3 Fix \mathbf{R}_{x} , \mathbf{O}_{x} , \mathbf{R}_{v} , \mathbf{O}_{v} , and \mathbf{S} , and update \mathbf{L} .

The algorithm is implemented in a coarse-to-fine manner, and two iterations are computed at each level of detail.

Results



only motion cues



reconstruction



A profile view of the The profile view with recovered albedo map



Reconstruction using both motion and photometric cues