

Final Project Proposal - Expanding Learning-Based Light Field View Synthesis

Thomas Lauer
UC San Diego
9500 Gilman Drive, San Diego CA
tlauer@ucsd.edu

Winston Durand
UC San Diego
9500 Gilman Drive, San Diego CA
wdurand@ucsd.edu

Abstract

Light field photography allows for the capture of images from multiple perspectives using a single camera array or array of microlenses. However these come at a tradeoff between spatial resolution and angular resolution. As the size of the microlens elements increases, the resulting light-field will have higher angular resolution which is useful for depth estimation, but will have lower spatial resolution for novel view synthesis. We will be reimplementing a paper [1] which uses a learning based method to synthesize views from light field camera data.

1. Introduction

A Lytro is a common light field camera which uses microlenses to capture multiple perspectives of a scene in one light field.

Many geometric methods such as Levoy *et al.* [2] require relatively high sample counts to ensure full coverage of the 4D lightfield, as the simple linear interpolation used by Levoy do not work well if there are heavily occluded regions or missing data.

Kalantari summarizes several existing methods for interpolating lightfields and their drawbacks. Many rely on high quality input images and well defined orientations between views, which are difficult to achieve with consumer light field cameras, and impossible with light fields captured by hand with cellphone cameras.

Wanner *et al.* [3] utilizes an optimization based approach, which calculates disparities using traditional computer vision as a preprocessing step. However, because the disparity estimation is independent from the loss function, it can not be optimized as part of the loss function.

Kalantari *et al.* [1] propose a method to interpolate between sparsely sampled sub-apertures views. They use an 8×8 subset of the full 14×14 subaperture, since the edge pixels are often black. The four corner sub-aperture views are fed into a series of two networks, the first is used to predict disparity which is then used to warp the sample im-

ages to the final perspective. The second network then takes these warped images, along with some additional metadata, and blends them together to produce the final RGB image of the novel view. Keeping this process differentiable allows both CNNs to be trained at the same time, which lets the disparity estimator become tuned to work with the final CNN.

2. Our Proposal

Our goal is to reimplement the paper *Learning-Based View Synthesis for Light Field Cameras* by Kalantari *et al.* [1] and additionally try to apply this technique to non-Lytro camera array captures. Further, we would like to investigate rendering novel views outside of the planar quadrilateral formed by the source images.

3. Milestones

Below are our milestones.

3.1. Set up original Matlab implementation (May 16th)

Our first goal will be to set up the original paper's code and recreate their results. This will give us a baseline to compare our new implementation to as we move forward. We will also study their implementation to see how they compute the gradient of the bicubic interpolation.

3.2. Naive CNN Approach (May 19th)

Our second goal will be to implement a single layer CNN similar to the Naive approach outlined in the paper. This is a single convolutional network which accepts the 4 sampled images along with their position information as input, and attempts to produce the novel view requested. As Kalantari *et al.* show, this should be able to create the novel view, but the result should be blurry and low detail. This is a good progress check to ensure we can actually train the network correctly.

We estimate that up to this point should take approximately a week.

3.3. Two-Part CNN With Disparity Estimation (May 25th)

The next step will be implementing the two-stage CNN architecture outlined by Kalantari *et al.* Instead of feeding the raw images from the sampled views, we apply a disparity warp to all the images at a set number of disparity levels. These are accumulated into a feature vector and used as input to the disparity prediction network. More details can be found in section 3.1 of Kalantari *et al.* [1].

3.4. Investigate Using Different Sampled Orientations (May 29th)

An additional approach we would like to consider to get more spread out of the Lytro capture data is to choose our input images in a diamond pattern rather than the 8×8 grid from the original paper.

3.5. Investigate Using Camera Array Datasets (May 29th)

As a followup, we intend to investigate how model trained in this manner is able to generalize to other camera array datasets, both where the views don't form a regular square and camera orientation varies.

3.6. Exploration of novel view outside capture bounds (Stretch goal)

Finally, we would like to experiment with changing the virtual camera's depth, moving it gradually away from the capture plane. A further path to test is the ability of our approach to estimate novel views which are still on the same plane as our reference views, but outside the bounds of the quadrilateral which they form. We could also test different numbers of sample images, the original paper used four to cover the entire 8×8 grid, but we could investigate how more or fewer images influence the quality of the results.

We generally are classifying this milestone as ambition.

4. Technology

We intend to implement this using PyTorch, because it has high performance, automatically handles calculating derivatives, and we have previous experience with it. Additionally, we will likely be using OpenCV or something similar for feature extraction in the preprocessing stage. We are not expecting to get real time performance, the paper claims it took 12.3 seconds to synthesize a single image using a Matlab implementation. We will stitch the individual images into a video to help visualize the results.

5. Datasets

Our training datasets will be Lytro captures since it represents a standardized input format which

is widely available, for example from Stanford at <http://lightfields.stanford.edu/LF2016.html>. Because

Lytro images are consistently formatted, we should be able to easily test our code with many different light fields. We also want to test this approach with a camera array, and we could easily adapt our preprocessing to accept the Stanford Multi-Camera Array. We will not be considering more free form light fields, such as those captured by handheld phone cameras, because of the added complexity of different camera orientations.

6. Questions to be Answered

How does the disparity from the first CNN compare to disparity produced by traditional computer vision methods, such as OpenCV?

Will the approach outlined work when we select sampled images from different locations, instead of the corners of the sub-aperture views? Further, is the CNN learning in this fixed grid transferrable to other grid layouts?

Will this work if we try to reconstruct a novel view outside the bounds of the sampled sub-aperture views?

If learning in non-transferrable across capture grid aspect ratios, is it possible to solve this with image stretching?

References

- [1] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM Transactions on Graphics (TOG)*, 35(6):1–10, 2016.
- [2] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 31–42, 1996.
- [3] Sven Wanner and Bastian Goldluecke. Variational light field analysis for disparity estimation and super-resolution. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(3):606–619, 2014.