

$$\sum_{t=1}^{T-1} \sum_p \left( \sum_{q \in \mathcal{N}_{p,\text{time}}^0} \phi_{\text{time}}^0(b_t^p, b_{t+1}^q) + \sum_{q \in \mathcal{N}_{p,\text{time}}^1} \phi_{\text{time}}^1(b_t^p, b_{t+1}^q) \right) + \sum_{t=1}^T \sum_p \left[ \sum_{q \in \mathcal{N}_{p,\text{space}}} \phi_{\text{space}}(b_t^p, b_t^q) + \phi_{\text{affine}}(b_t^p) \right]. \quad (7)$$

that the configuration is unchanged when  $b_t^p = 0$ , and an alternative layer  $\hat{k}$  becomes visible when  $b_t^p = 1$ . Revealing this new layer may alter all the support functions for the first  $\hat{k}$  layers: when  $b_t^p = 1$ ,  $g_{t\hat{k}}^p(1) = 1$  and  $g_{t\hat{k}}^p(1) = 0, k < \hat{k}$ . In this case, we also set the motion for layer  $\hat{k}$  to that of the formerly visible layer ( $u_{t\hat{k}}^p(1) = u_{t\hat{k}'}^p(1)$ ), and the motion for layer  $k'$  to its affine mean ( $u_{t\hat{k}'}^p(1) = u_{\theta_{t\hat{k}'}}^p(1)$ ).

This segmentation and flow move involves many terms from the overall model, which depend on the binary decision variables as summarized in Eq. (7). The choice of  $b_t^p$  influences the flow vector at pixel  $p$ , and thus determines which of two candidate pixels it is linked to at the next frame. When  $b_t^p = 0$ , the temporal neighbors are

$$\mathcal{N}_{p,\text{time}}^0 = \{(i + [u_{t\hat{k}}^p(0)], j + [v_{t\hat{k}}^p(0)]), 1 \leq k \leq K\}, \quad (8)$$

and the corresponding potential function will only “fire” when  $b_t^p = 0$ , so that  $\phi_{\text{time}}^0(b_t^p, b_{t+1}^q) =$

$$\left[ \left( \rho_d(\mathbf{I}_t^p - \mathbf{I}_{t+1}^{q'}) - \lambda_d \right) s_{t\hat{k}}^p(b_t^p) s_{t+1,\hat{k}}^q(b_{t+1}^q) + \lambda_c (1 - \delta(g_{t\hat{k}}^p(b_t^p), g_{t+1,\hat{k}}^q(b_{t+1}^q))) (1 - \delta(k, K)) \right] (1 - b_t^p), \quad (9)$$

which incorporates both the data and temporal terms for the  $K - 1$  support functions. We evaluate the warped image  $I_{t+1}^{q'}$  at subpixel positions and the visibility mask  $s_{t+1,\hat{k}}^q$  and the warped support function  $g_{t+1,\hat{k}}^q$  at integer positions. Similarly  $\phi_{\text{time}}^1(b_t^p, b_{t+1}^q)$  is defined to only “fire” when  $b_t^p = 1$  (see **Supplemental Material**).

For the spatial term, the set  $\mathcal{N}_{p,\text{space}}$  contains the four nearest neighbors of pixel  $p$ . The binary selection variable changes the states of several binary support functions and flow fields. The effects sum together as

$$\phi_{\text{space}}(b_t^p, b_t^q) = \sum_{k=1}^{K-1} \lambda_b w_q^p (1 - \delta(g_{t\hat{k}}^p(b_t^p), g_{t\hat{k}}^q(b_t^q))) + \sum_{k=1}^K \lambda_a \rho_{\text{mrf}}(u_{t\hat{k}}^p(b_t^p) - u_{t\hat{k}}^q(b_t^q)). \quad (10)$$

The unary term can be obtained from Eq. (6) as

$$\phi_{\text{affine}}(b_t^p) = \sum_{k=1}^K \lambda_a \lambda_{\text{aff}} \rho_{\text{aff}}(u_{t\hat{k}}^p(b_t^p) - u_{\theta_{t\hat{k}}}^p). \quad (11)$$

**Visibility move.** Given the current flow estimate, we decide whether to make a pixel  $p$  visible for the selected layer  $\hat{k}$  by modifying the previous layer support  $\mathbf{g}^{\text{old}}$ . When  $b_t^p = 0$ , all the support functions retain their previous value at  $p$ ,

*i.e.*  $g_{t\hat{k}}^p(0) = g_{t\hat{k}}^{p,\text{old}}$ . When  $b_t^p = 1$ , we need to adjust the support functions of the first  $\hat{k}$  layers so that layer  $\hat{k}$  is visible at  $p$ , *i.e.*  $g_{t\hat{k}}^p(1) = 0$  if  $k < \hat{k}$ , and  $g_{t\hat{k}}^p(1) = 1$  if  $k = \hat{k}$ . When  $\hat{k}$  is the last layer, all the support functions of the first  $K - 1$  layers are set to be 0 at  $p$ .

**Flow selection move and continuous refinement.** Given the current layer assignment, the pixel of each layer can retain its current motion or take the motion of the affine mean flow field. This is a binary segmentation problem similar to the FusionFlow [17] work. The main difference is that our model uses the segmentation information to handle occlusions. After this step, we refine the output flow field by continuous optimization [26] to adaptively change the candidate flow fields for the discrete optimization.

### 3.3. Layer Number Determination and Depth Order Reasoning

We initialize with an upper bound on the number of layers. During optimization, when a layer has no visible pixels associated with it, we remove it from the solution. The new solution can equally explain the image data, pays no penalty for the removed layer, and so has lower energy. Inferring the depth ordering of layers requires testing all the possible combinations and is computationally prohibitive. We instead use heuristics to reduce the search space. We first order the layers from fast to slow by their average motion. We then perform the moves above to estimate the support functions and the flow fields in both the fast-to-slow and the slow-to-fast ordering. The ordering with the lower energy is further refined as follows. For each pair of neighboring layers, we propose to switch their ordering, and optimize their visibility mask and support functions. If the new solution has a lower energy than its previous one, we accept this new depth ordering and proceed to other pairs (see **Supplemental Material** for the detailed algorithm). In practice, we find that this local greedy search scheme is fairly robust.

## 4. Experimental Results

We evaluate the proposed layered model on both motion estimation and layer segmentation tasks. Throughout this section, the proposed method is called **nLayers**, since it can automatically determine the number of layers. **Layers++** refers to the continuous method developed in [27] which uses a fixed number of 3 layers. For layer segmentation, we also compare our method to a state-of-the-art, hierarchical graph-based video segmentation algorithm [9], referred to