# Neural Scene Flow Fields for Space-Time View Synthesis of Dynamic Scenes

## Related Work

view synthesis   nerf ( static) $\begin{cases} \text{kernels btw two images} \\ \text{optical flow, warping features} \\ ? \text{ layered} \end{cases}$

novel time synthesis

space-time view synthesis   scene change with time
$\begin{cases} \text{illumination changes: previous relighting work only assume static} \\ \text{3D scene motion: previous require multi-view, time synchronised as input} \end{cases}$

## Approach

### Neural scene flow fields for dynamic scenes (nsff. the dynamic version of nerf)

recall: $(c, \sigma) = F_\Theta(x, d)$   # nerf takes a position, a direction. outputs a rgb, a density.

now: $(c_i, \sigma_i, F_i, W_i) = F_\Theta^{dyn}(x, d, i)$   # nsff $F_\Theta^{dyn}$ takes additionally a time $i$, outputs additionally:

$\begin{cases} F_i = (f_{i \to i+1}, f_{i \to i-1}) & \text{\# 3D scene flow, offset of location } X @ i-1, i+1 \\ W_i = (w_{i \to i+1}, w_{i \to i-1}) & \text{\# 1D disocclusion weights } W_i. \text{ details below} \end{cases}$

## Optimization   core of nsff. to understand this paper from THE NEW LOSS

### temporal photometric consistency   scene @ $i$ should be consistent with scene @ $j \in N(i)$, if neighbor scene got from nsff.

$\hat{C}_{j \to i}(r_i) = \int_{t_n}^{t_f} T_j(t) \sigma_j(r_{i \to j}(t)) c_j(r_{i \to j}(t), d_j) dt$   where $r_{i \to j}(t) = r_i(t) + f_{i \to j}(r_i(t))$   noticing $t$ is distance (consider as a position is good)   # color for $r_i$ is from time neighbors

roughly: $L_{pho} = \sum_{r_i} \sum_{j \in N(i)} \| \hat{C}_{j \to i}(r_i) - C_i(r_i) \|_2^2$   # ray-wise, all neighbors should render same as ground truth for a single pixel

However. this is not always correct for neighboring,   # motion causes 3D disocclusion regions, this is ambiguity. Or to say:

So, results from neighbors should be weighted (not always average)

> you cannot ask neighbors render exact same! what if there is occlusion/disocclusion?

Weights for each neighbor: $\hat{W}_{j \to i}(r_i) = \int_{t_n}^{t_f} T_i(t) \sigma_j(r_{i \to j}(t)) W_{i \to j}(r_i(t)) dt$
   $\overbrace{[0,1] \text{ to } 1}$

# you should be noticing rgb is function of direction. but this weight is not

# for all pixels, your render results should be weighted neighbor-consistent.

# for all rays, your render points should have everywhere close-to-1 weights

the final: $\boxed{L_{pho} = \sum_{r_i} \sum_{j \in N} \hat{W}_{j \to i}(r_i) \| \hat{C}_{j \to i}(r_i) - C_i(r_i) \|_2^2 + \beta_w \sum_{x_i} \| w_{i \to j}(x_i) - 1 \|_1}$

for all pixels each ray      for all rays each position

# $\beta_w$: reg weight (= 0.1). $N(i) = 5$. ($i \pm 2$ chained?).

# $j = i$. self case, $\hat{W}_{i \to i}(r_i) = 1$, $f_{i \to i} = 0$. $\hat{C}_{i \to i}(r_i) = \hat{C}_i(r_i)$. original case.

### Scene flow priors   points back & forth with nsff should be consistent with each other

$x_i$ have forward scene flow $f_{i \to j}$. gives point $x_{i \to j} = x_i + f_{i \to j}$

$x_i \underset{f_{j \to i}}{\overset{f_{i \to j}}{\rightleftarrows}} x_j (= x_{i \to j})$   this has to be true

$f_{i \to j}(x_i) + \underline{f_{j \to i}(x_j)} = f_{i \to j}(x_i) + f_{j \to i}(x_{i \to j}) = 0$

   no need to really bother $j$ now

> you are predicting next x position with forward flow, then take the new next position. Naturally, should the new position take one step back at exactly x.

$\boxed{L_{cyc} = \sum_{x_i} \sum_{j, i \neq j} w_{i \to j} \| f_{i \to j}(x_i) + f_{j \to i}(x_{i \to j}) \|_1}$   # for each position, this should be true always. only closest neighbors 2

## data-driven priors   results better consistent with other methods w.r.t. positions

### geometric consistency   2D optical flow guides 3D scene flow fields

> optical flow methods give a next-time 2D position, nsff gives a 3D position but we do projection. They should be same.

2D optical flow   $P_{i \to j} = P_i + u_{i \to j}$ ( Given By OF methods)

3D nsff   for one ray $r_i$. we volume render the $f_{i \to j}$ and $x_i$   # recall how we compute depth? recall volume render seems apply everything

$\hat{F}_{i \to j}(r_i) = \int_{t_n}^{t_f} T_i(t) \sigma_i(r_i(t)) f_{i \to j}(r_i(t)) dt$   next world position: $\hat{X}_i(r_i) + \hat{F}_{i \to j}(r_i)$, then, do projection to frame at time $j$.

$\hat{X}_i(r_i) = \int_{t_n}^{t_f} T_i(t) \sigma_i(r_i(t)) x_i(r_i(t)) dt$

$\hat{P}_{i \to j}(r_i) = \pi \left( K \left( R^j_c (\hat{X}_i(r_i) + \hat{F}_{i \to j}(r_i)) + t^j \right) \right)$   $\begin{cases} [R^j, t^j] \text{. extrinsic at time } j \\ K \text{. intrinsic shared} \end{cases}$   # world position (trained from nerf/nsff) to camera position.

$\boxed{L_{geo} = \sum_{r_i} \sum_{j \in N_2} \| \hat{P}_{i \to j}(r_i) - P_{i \to j}(r_i) \|_1}$   # for all pixels/each ray result, it should be same.

### single-view depth   monodepth prediction guides

$\boxed{L_z = \sum_{r_i} \| \hat{z}_i^*(r_i) - \hat{z}_i^*(r_i) \|_1}$   # check the supp for real implementation here.

## Integrating a Static Scene Representation — differentiate background static and foreground dynamic then combine

Static background scene: $(c, \sigma, v) = F_\theta^{st}(x, d)$  <span style="color:red"># static/dyn balance : $v$  (fyi, this is not a necessarily trained param)</span>

dynamic foreground scene: $(c_i, \sigma_i, F_i, w_i) = F_\theta^{dy}(x, d, i)$  <span style="color:red"># blending weight $v$ is actually $v = f(c, \sigma_i)$</span>

Combined rendering equation: $\hat{C}_i^{cb}(r_i) = \int_{t_n}^{t_f} T_i^{cb}(t) \, \underline{\sigma_i^{cb}(t) \, c_i^{cb}(t)} \, dt.$  ↗ a weighted render

$$L_{cb} = \sum_{r_i} \| \hat{C}_i^{cb}(r_i) - C_i(r_i) \|_2^2$$

<span style="color:red"># this is a PURE color render, basically has a PURE nerf for background. and a nsff for foreground.
we don't render with neighbor information!</span>

## Regularisation

### spatial smoothness — points along ray should have neighboring similar scene flow

$$L_{sp} = \sum_{r_i} \sum_{y_i \in N(x_i)} \sum_{j \in i \pm 1} \exp(-2\|x_i, y_i\|_2) \| f_{i \to j}(x_i) - f_{i \to j}(y_i) \|_1$$

<span style="color:red">dist</span>  <span style="color:red"># for each ray, considering each point</span>

### temporal smoothness — backward and forward should equal <span style="color:red"># for some reason</span>

$$L_{temp} = \frac{1}{2} \sum_{x_i} \| f_{i \to i+1}(x_i) + f_{i \to i-1}(x_i) \|_2^2$$

<span style="color:red"># for each point</span>

### sf minimal — encourage sf to be low <span style="color:red"># cause frames are continuous</span>

$$L_{min} = \sum_{x_i} \sum_{j \in i \pm 1} \| f_{i \to j}(x_i) \|_1$$

<div style="border:1px solid red">

**FINALLY**

$$L = L_{cb} + L_{pho} + \beta' L_{cyc} + \beta'' L_{data} + \beta''' L_{reg}$$

where $\begin{cases} L_{data} = L_{geo} + \beta_z L_z \\ L_{reg} = L_{sp} + L_{temp} + L_{min} \end{cases}$

</div>

## rendering — the splatting-based plane-sweep volume rendering approach <span style="color:red"># render with changing time. not view</span>

to render intermediate time $i + \delta i$. $\delta i \in (0,1)$, sweep over every step $t$ (a location along the ray)

for one point $\begin{cases} \text{query } i \text{ and get}: (c_i, \sigma_i) \, (F_i) \\ \text{query } i+1 \text{ and get}: (c_{i+1}, \sigma_{i+1})(F_{i+1}) \end{cases}$

at time $i + \delta i$, points will be <u>NEW</u> and <u>BUFFERED</u>.

$(c, \sigma)$ given by queried results already. $\begin{cases} x_i + \delta i f_{i \to i+1}(x_i) \\ \text{position given by new as,} \\ x_{i+1} + (1-\delta i) f_{i+1 \to i} \end{cases}$

<u>two parts will be linearly blended</u>. <span style="color:red"># for a buffer</span>

## Implementation Details

refer to supp.

<div style="border:1px solid red; color:red">

How to understand this algo?

basically we don't know the color @ $i + \delta i$ and how color flows, so we query
@ $i$ and @ $i+1$. the result includes $(c, \sigma)$ and $(f, w)$. so now, $f$ is
indicating where these points in nsff are going to. for time $i + \delta i$, the
querying results come from a BUFFER. this BUFFER has the points
warped by $f$ from time $i$ and $i+1$, along with their original queried $(c, \sigma)$.
By this way, we don't query directly from nsff for $(c, \sigma)$ at intermediate
time. we only query a correct 'NERF' result and correct 'NSFF' flow to
get what we want finally

</div>