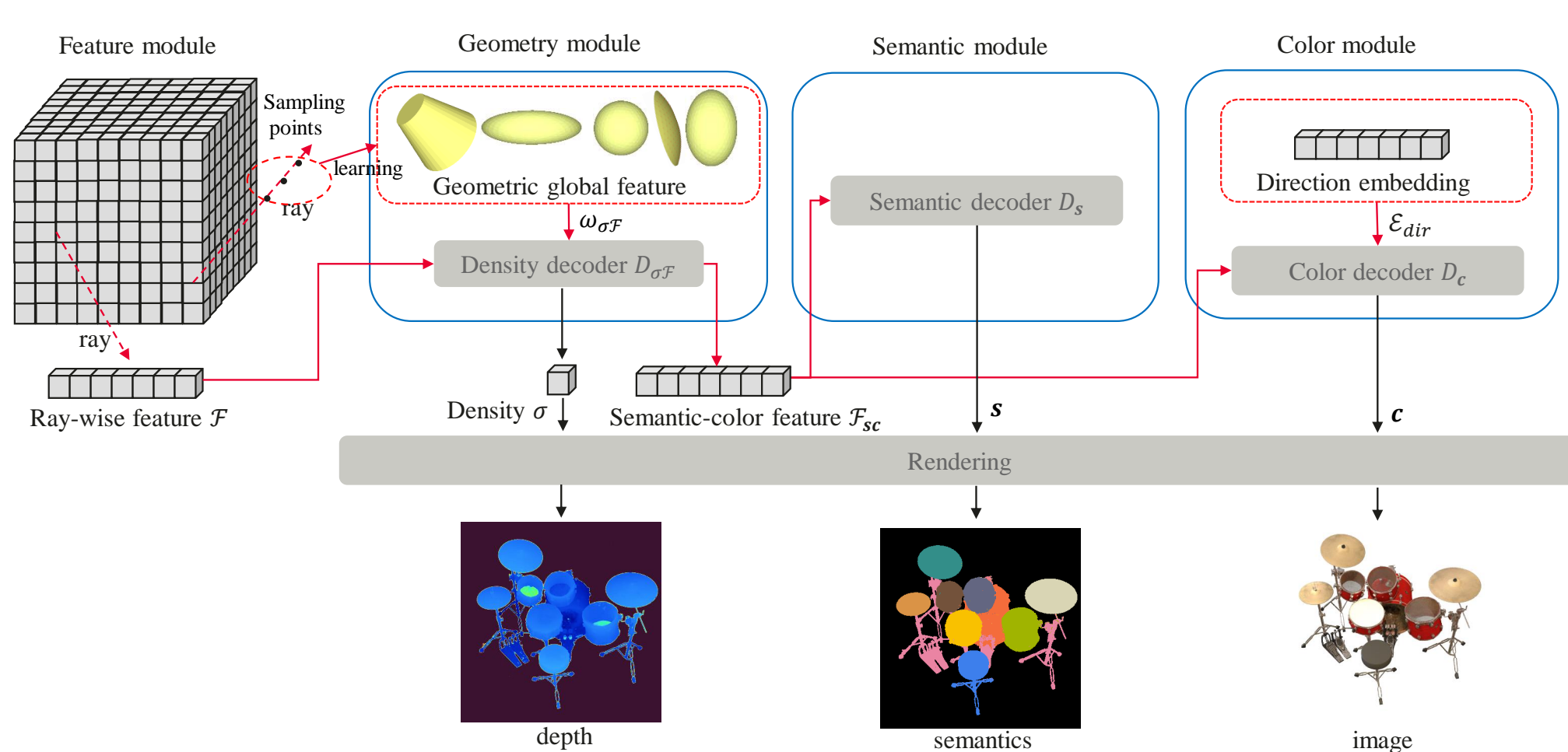


## IMPLICIT DATA ENGINE



The contributions of this paper:

- IS-NEAR is the first to associate all independent NeRF rays with 3D global features, which improves the geometry, color, semantics and labeling efficiency.
- We customize the back-propagation to eliminate the differences between geometric and semantic inferences.
- The carefully loss and network design makes a trade-off among efficiency, semantics, color and geometry, superior to the SOTA methods.
- The engine can be applied to indoor, outdoor and objects semantic labeling, texture re-rendering, and robot simulation.

## 3D GLOBAL FEATURES

The point-to-surface representation is expressed as:

$$d = \pi^T \cdot \mathbf{X}, \quad (1)$$

where  $\pi$  is the coefficients vector of the quadratic terms  $\mathbf{X}$ , namely the global feature,  $\mathbf{X}$  is expressed as:

$$\mathbf{X} = (x^2, y^2, z^2, xy, xz, yz, x, y, z), \quad (2)$$

$x, y, z$  is the sampling points on each ray.

$$\omega = L_{emb}(1 - (\text{sigmoid}(\pi^T \cdot \mathbf{X}))). \quad (3)$$

$\omega$  is related to the distance between the point and the global surface.

## IMPLICIT FIELDS

The point density  $\sigma$ , semantics  $s$  and color  $c$  are defined as:

$$[\sigma, \mathcal{F}_{sc}] = \omega_{\sigma\mathcal{F}} \cdot D_{\sigma\mathcal{F}}(\mathcal{F}), \quad (4)$$

$$s = D_s(\mathcal{F}_{sc}), \quad (5)$$

$$c = D_c([\mathcal{F}_{sc}, \mathcal{E}_{dir}]), \quad (6)$$

## LOSSES

$$\mathcal{L}_c = \frac{1}{B} \sum_{b=1}^B \|c_b - \hat{c}_b\|^2, \quad (7)$$

$$\mathcal{L}_s = - \sum_{r \in \mathcal{R}} \sum_{k=1}^{N_c} w_k \hat{p}_k(r) \log p_k(r), \quad (8)$$

$$w_k = \lfloor \frac{n_k}{\sum_{i=0}^{N_c} n_i} \rfloor^h, \quad (9)$$

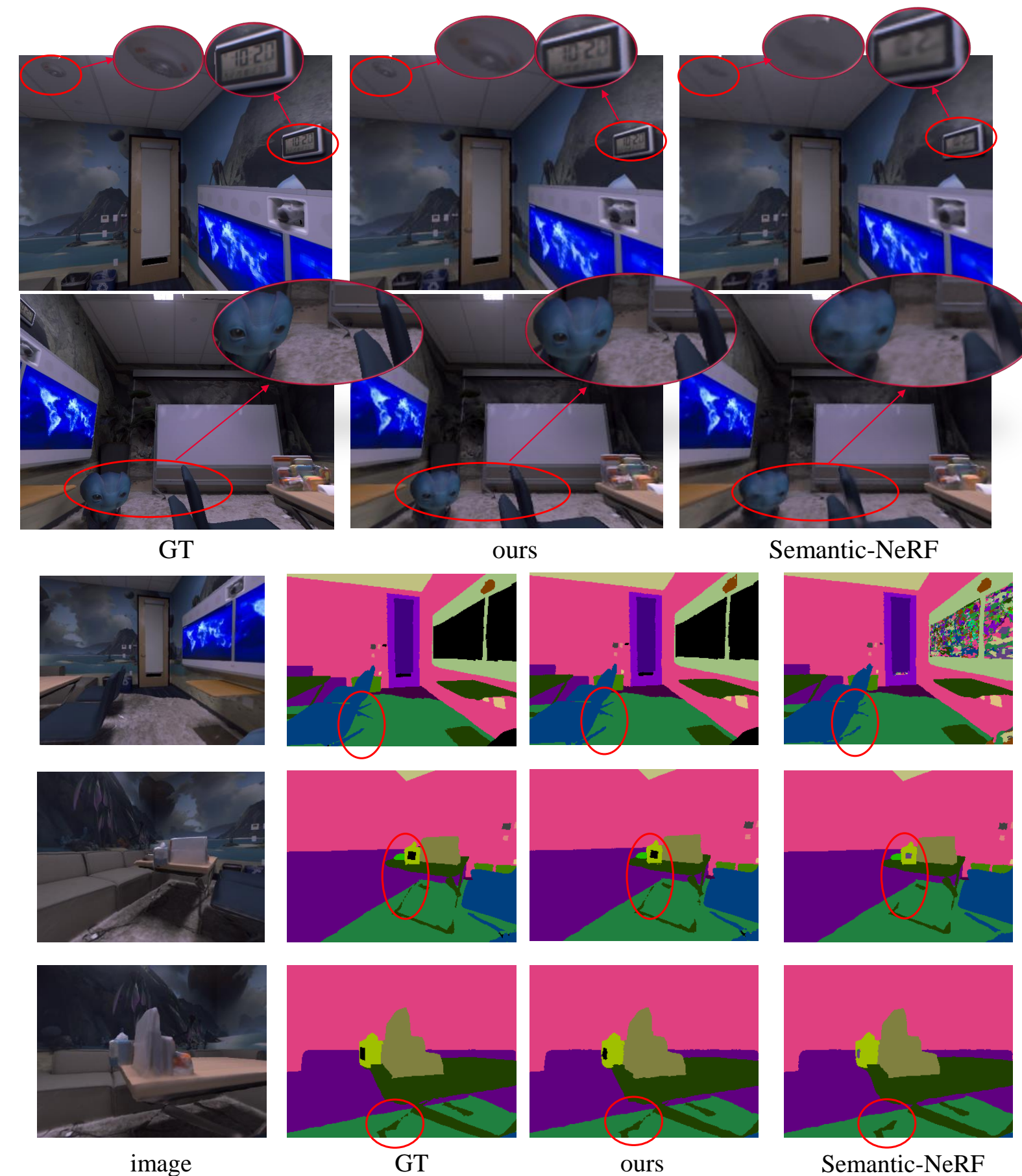
$$\lfloor \lfloor x \rfloor \rfloor^h = \begin{cases} l, & x < l \\ x, & l \leq x \leq h \\ h, & x > h \end{cases}. \quad (10)$$

$$\mathcal{L}_d = \sqrt{\frac{1}{N_d} \sum_i g_i^2 + \frac{\lambda}{N_d^2} (\sum_i g_i)^2}, \quad (11)$$

## COMPARISON RESULT

Performance Method	Efficiency		Color		Semantics		
	$T_t$	$T_r$	PSNR (dB)	SSIM	mIoU (%)	$ACC_t$ (%)	$ACC_a$ (%)
SS-NeRF	9 h	—	30.2	—	92.4	—	—
Semantic-NeRF	8 h	4.82 s	31.39	0.930	93.68	99.00	96.53
Ours	15 min	0.1 s	35.9	0.970	94.79	99.29	97.68

Performance Method	Geometry		
	AbsDiff	SqRel	RMSE
SS-NeRF	—	—	—
Semantic-NeRF	0.032	0.007	0.096
Ours	0.0056	0.0006	0.0122



## ROBUSTNESS TO SPARSE LABELS

The proposed IS-NEAR with point-to-surface global feature is robust to the sparse labels.

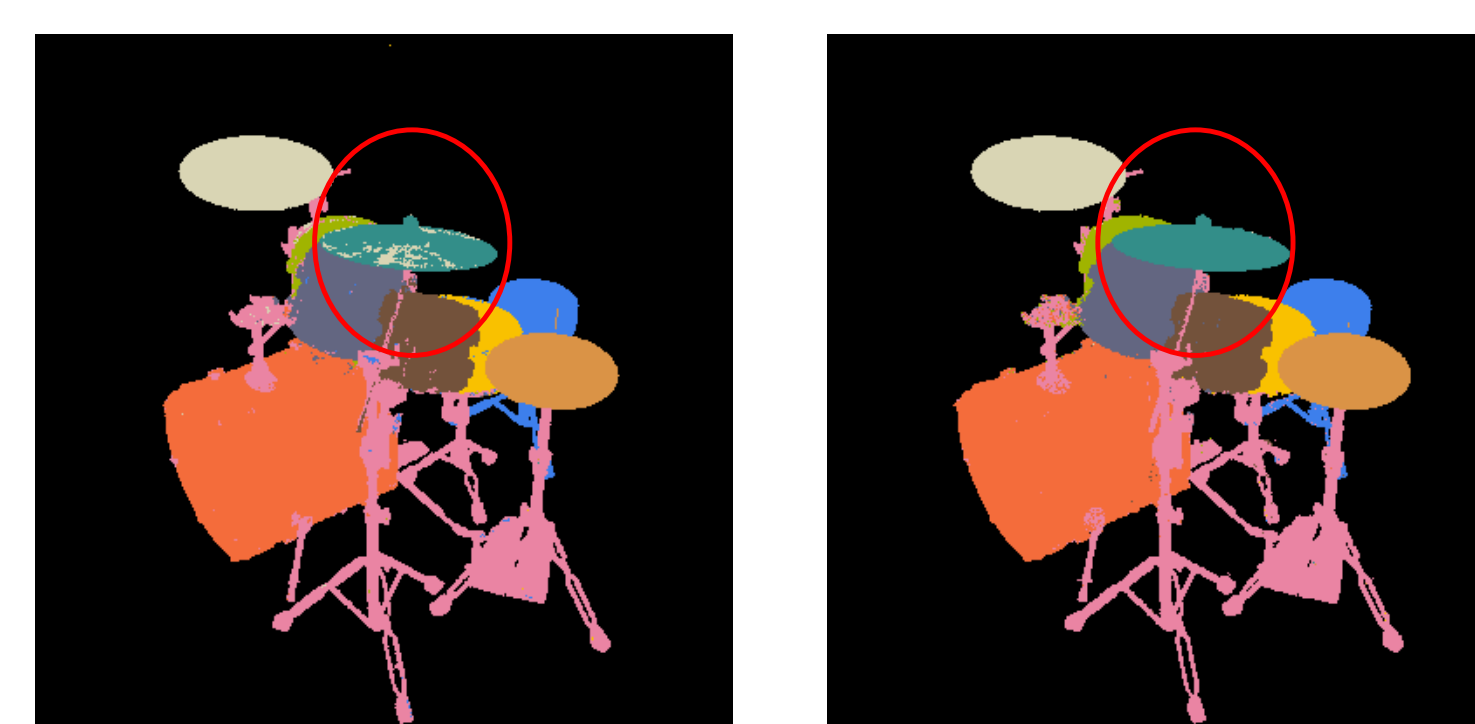
R	PSNR (dB)	SSIM	mIoU (%)	$ACC_t$ (%)	$ACC_a$ (%)
10%	40.7	0.989	93.1	99.1	96.1
20%	40.5	0.988	93.1	99.1	96.2
25%	40.3	0.988	92.4	99.0	96.2
50%	40.7	0.989	92.6	99.0	96.1
100%	40.6	0.989	92.3	99.0	95.9
10%w/o PS	39.6	0.985	83.6	97.3	91.0
20%w/o PS	39.1	0.983	89.3	98.5	94.5
25%w/o PS	39.1	0.982	89.4	98.6	94.6

## TRUNCATED-WEIGHTS SEMANTIC LOSS



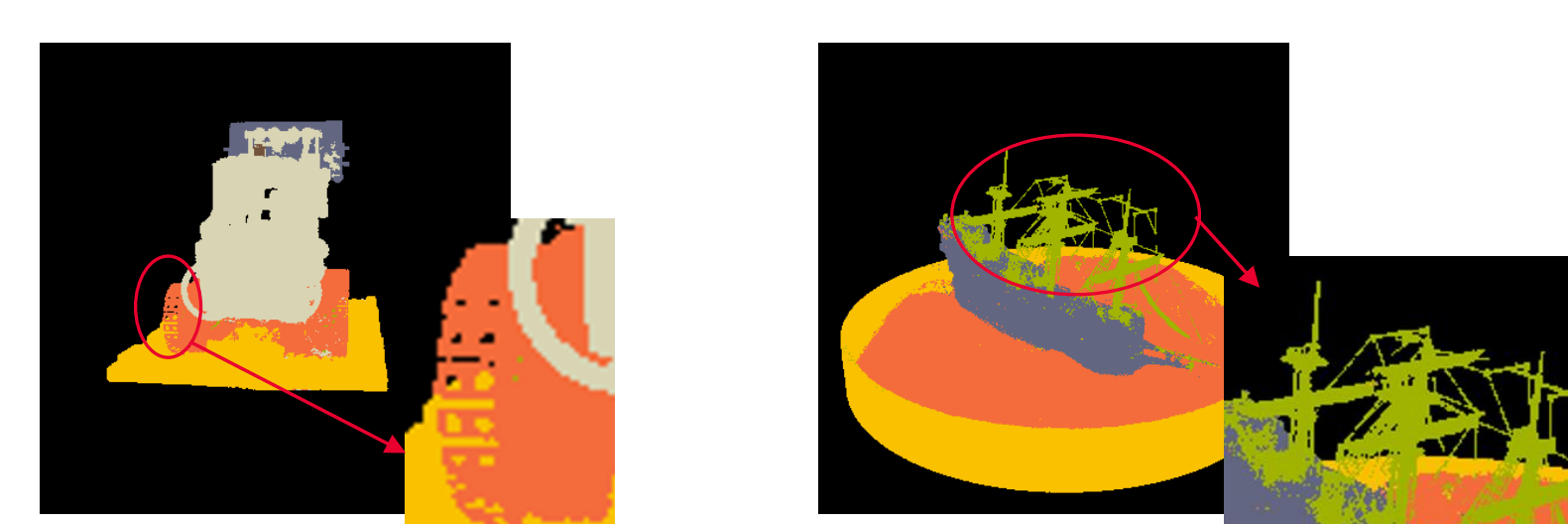
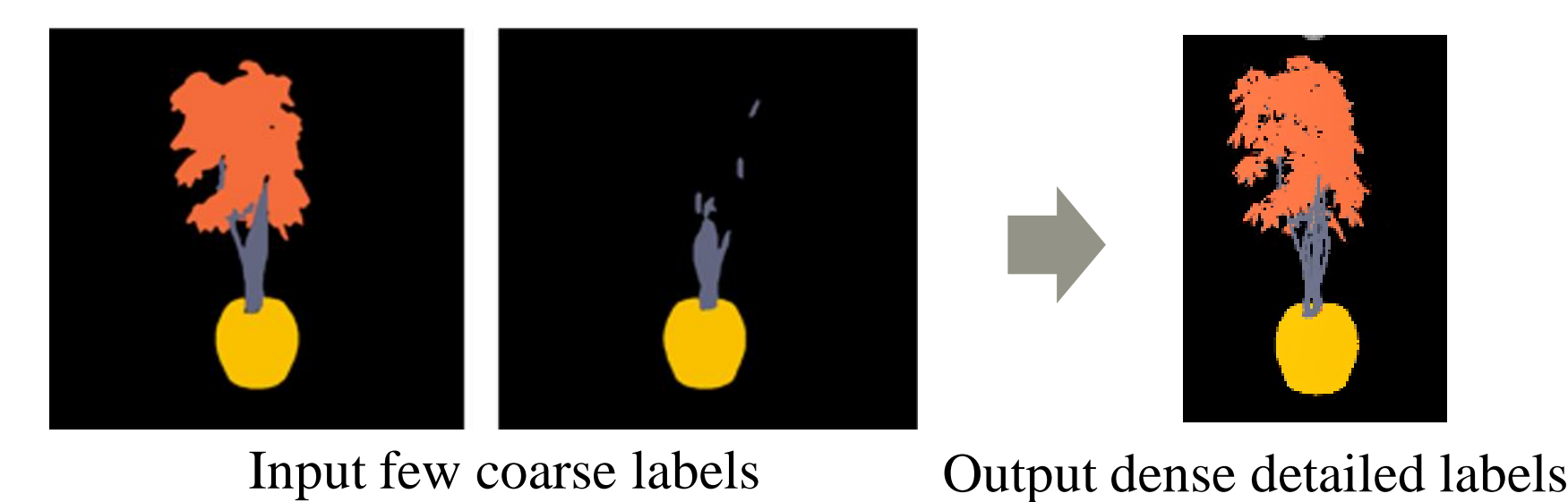
(a) Traditional cross-entropy loss  
(b) Truncated-weights cross-entropy loss (Eq. 10)

## SEMANTIC RENDERING WITHOUT DIRECTION EMBEDDING



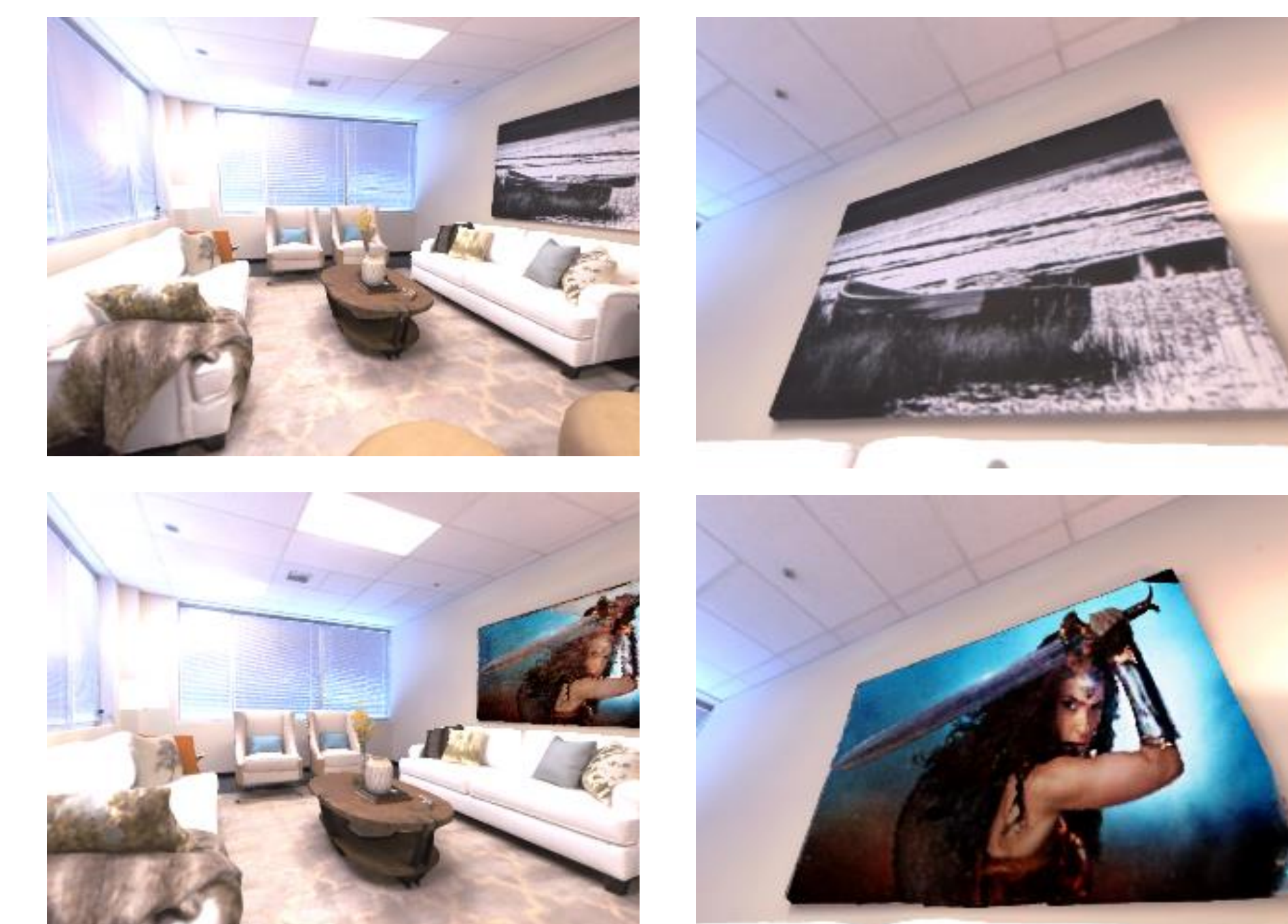
(a) With direction embedding (b) Without direction embedding

## APPLICATION-LABELING



Rendering detailed semantics

## APPLICATION-TEXTURE RERENDERING



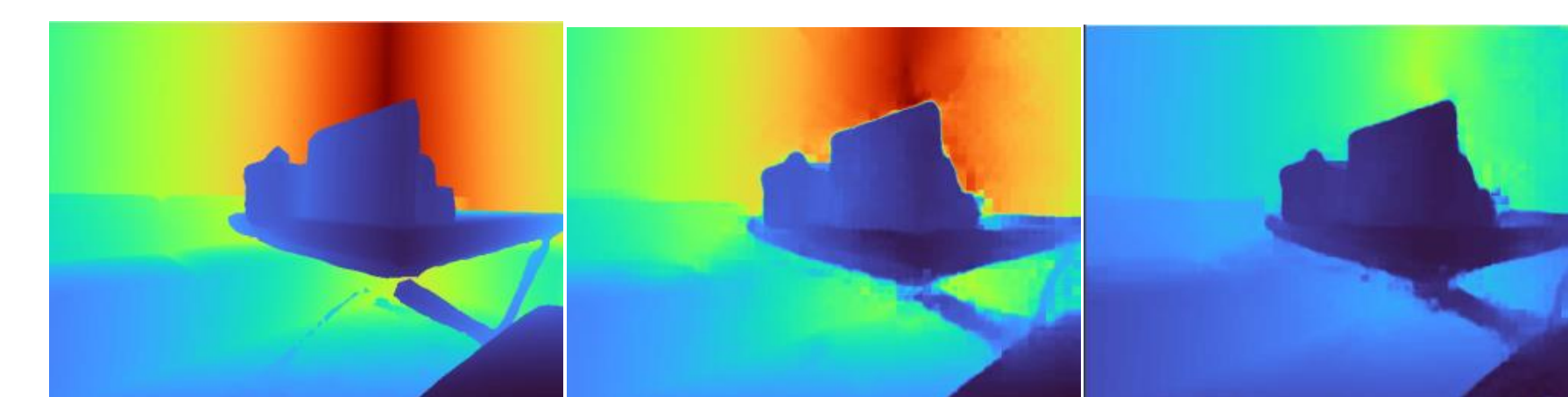
Update specific semantic region textures

## APPLICATION-SIMULATION

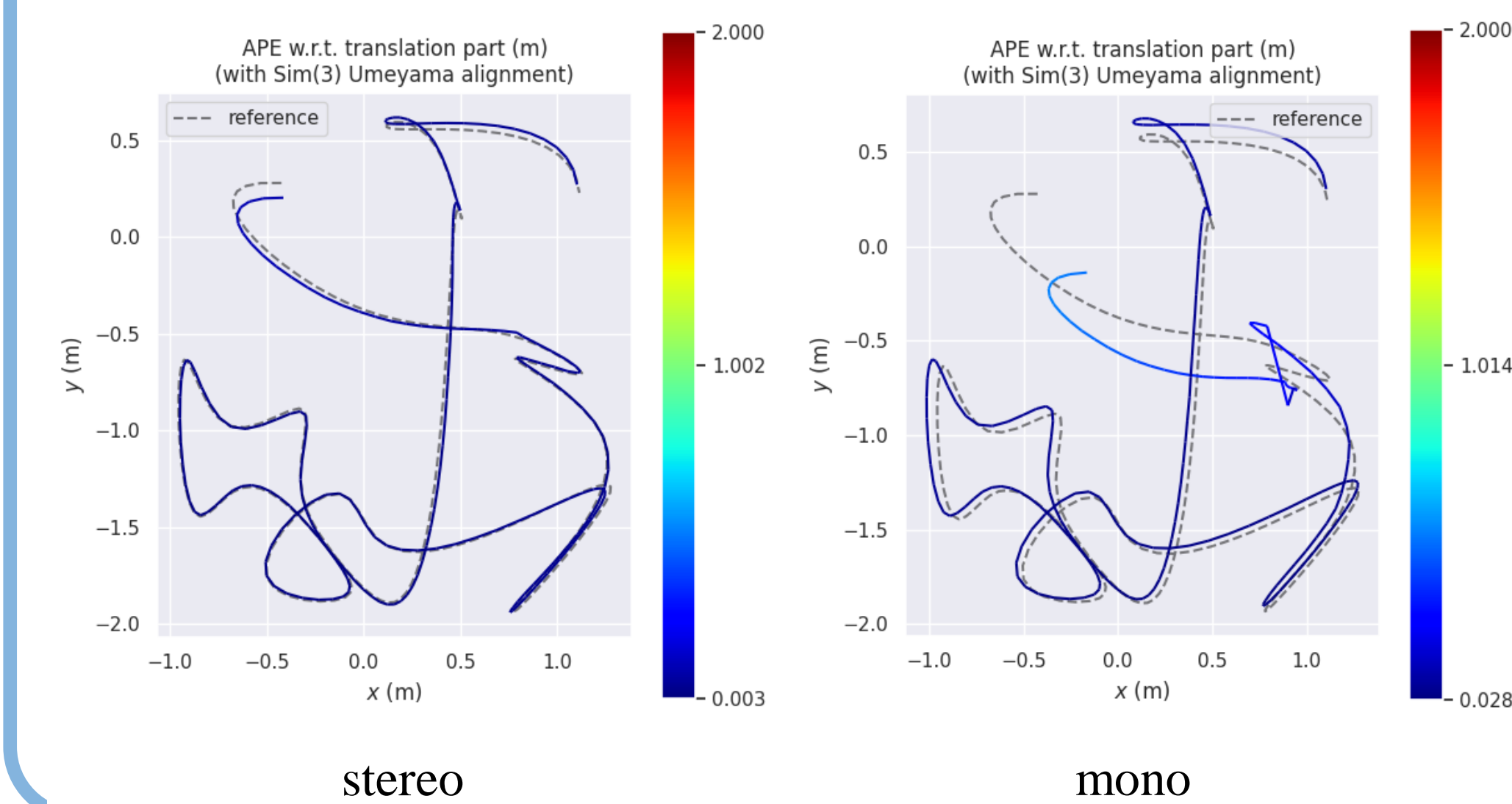
The monocular and binocular sequence data generated by IS-NEAR are used for SLAM simulation.



(a) (b)

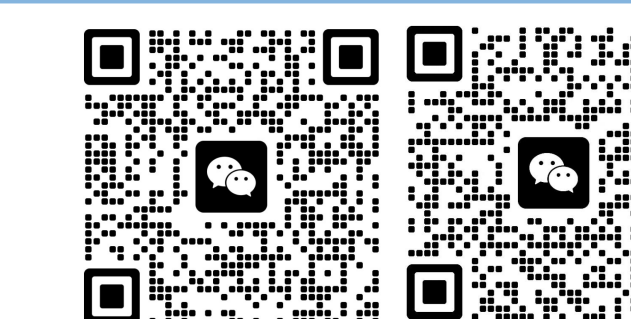


(c) (d) (e)



stereo mono

## FOR PAPER, RESULTS, CODE AND MORE



{suntiecheng1, lintao879, dongxingliang}@huawei.com  
wei.zhang@ifp.uni-stuttgart.de