

Optimal Transport for Rectified Flow Image Editing

A Unified Framework for Inversion-Based and Inversion-Free Editing

Marian Lupascu and Mihai-Sorin Stupariu

WACV 2026 — Paper #1419

github.com/marianlupascu/OT-RF

Why is Rectified Flow Editing Hard?

Rectified Flow models (FLUX, SD3) use straight-line trajectories — great for generation, but editing is challenging:

Problem 1 — Inversion-based

Deterministic ODE \Rightarrow no stochastic regularization.
Errors compound \Rightarrow **trajectory deviation**.

Problem 2 — Inversion-free (FlowEdit)

Pathways constructed *without* principled guidance.
 \Rightarrow **Suboptimal paths & semantic drift**.

Core trade-off: Reconstruction fidelity vs. Editing flexibility.



Key Insight: OT theory gives principled, minimal-cost pathways — for *both* editing paradigms, *training-free*.

Transport-Guided Inversion

OT correction at each denoising step:

$$v_{\text{enh}} = v_{\text{RF}} + \alpha(t) \cdot d_{\text{OT}}$$

Adaptive cosine schedule $\alpha(t)$. No retraining.
 $O(n)$ cost.

Transport-Enhanced FlowEdit

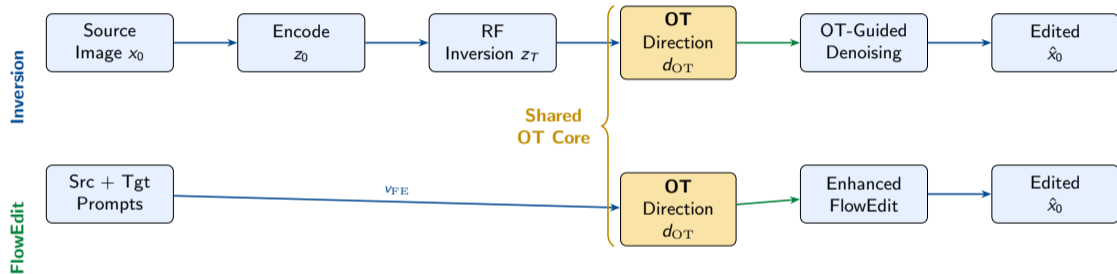
OT augmentation of FlowEdit velocity:

$$v_{\text{enh}} = v_{\text{FE}} + \gamma(t) \cdot d_{\text{OT}}$$

Better semantic consistency. Training-free.

Closed-Form Transport Direction (Brenier map, $W_2, O(n^3) \rightarrow O(n)$)

$$d_{\text{OT}}(z_t, z_{\text{tgt}}, t) = \frac{z_{\text{tgt}} - z_t}{\max(T - t, \delta)}$$



One transport core, two paradigms. OT guidance acts as *elastic constraints* keeping edits on the source manifold.

Reconstruction Quality — SFHQ Dataset (null prompt)

Method	LPIPS↓	SSIM↑	PSNR↑	Face Rec.↓	CLIP-I↑	RT (s)↓
RF-Inversion	0.135	0.833	32.58	0.387	0.936	28.8
RF-Inv.+OTC (ours)	0.001	0.992	40.71	0.112	0.999	28.6
FlowEdit (SD3)	0.001	0.990	40.78	0.094	0.998	16.5
FlowEdit (SD3)+OTC (ours)	0.001	0.990	40.78	0.093	0.998	16.6
FlowEdit (FLUX)	0.096	0.879	24.26	0.304	0.868	28.6
FlowEdit (FLUX)+OTC (ours)	0.001	0.993	40.65	0.113	0.999	28.7

RF-Inversion+OTC

×135 LPIPS +8.1 dB PSNR -71% face dist.

FlowEdit (FLUX)+OTC

×96 LPIPS +16.4 dB PSNR -63% face dist.

Stroke-to-Image Reconstruction

(LSUN-Bedroom / Church, L2 ↓)

Method	Bedroom	Church
RF-Inversion	82.55	80.35
+OTC	76.10	69.97
<i>gain</i>	-7.8%	-12.9%
FlowEdit (FLUX)	61.51	60.76
+OTC	57.61	57.43

Semantic Face Editing

Prompt: "...wearing glasses"

Method	Face Rec.↓	CLIP-T↑
RF-Inversion	0.621	0.281
+OTC	0.552	0.285
<i>gain</i>	-11.2%	—
FlowEdit (FLUX)	0.537	0.293
+OTC	0.521	0.296

- ✓ Better identity preservation
- ✓ Stronger text-image alignment

Qualitative Results — Reconstruction & Stroke-to-Image

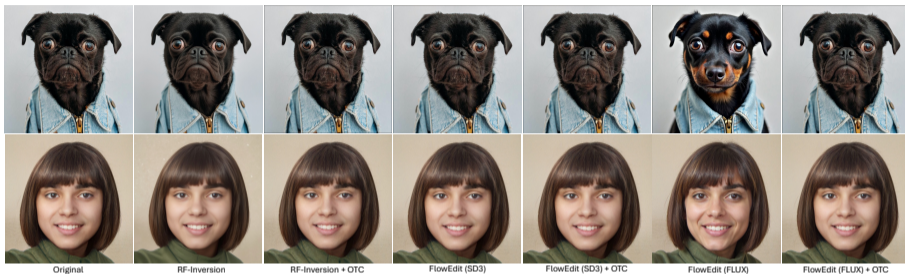


Fig. 3: Reconstruction across inversion-based (RF-Inversion, +OTC) and inversion-free (FlowEdit SD3/FLUX \pm OTC) methods.

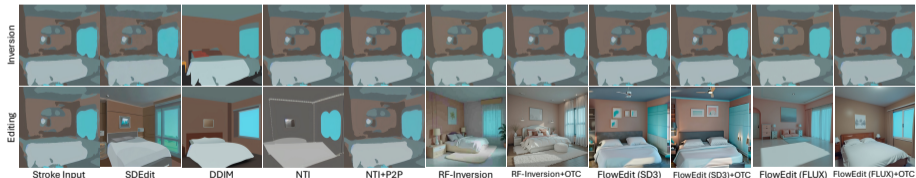


Fig. 4: Stroke-to-image: null prompt (top) vs. "a photo-realistic picture of a bedroom" (bottom).

✓ OTC consistently restores fidelity across **both inversion and FlowEdit** pipelines

Qualitative Results — Semantic Face Editing

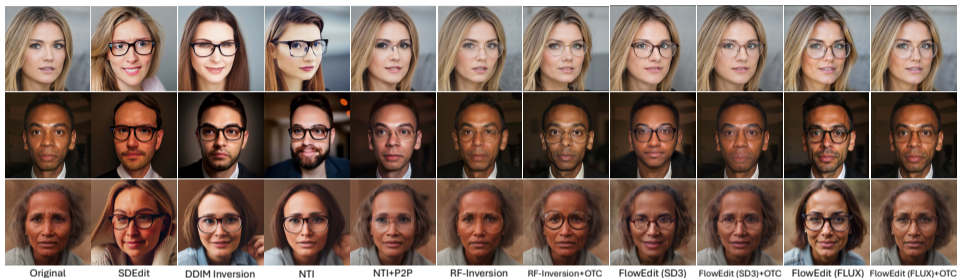


Fig. 6: Prompt “face of a man/woman wearing glasses” — identity preserved across 3 subjects, all methods.



Fig. 7: RF-Inversion (top) vs. RF-Inversion+OTC (bottom): expression, age, gender, object insertion.

- ✓ Better identity preservation
- ✓ Expression, age, gender, objects — all improved with OTC

Closed-Form W_2 Transport

Brenier map: optimal displacement is

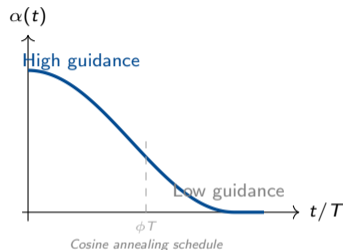
$$d_{OT} = (z_{tgt} - z_t) / (T - t)$$

- ✓ No iterative solver
- ✓ $O(n^3) \rightarrow O(n)$ complexity
- ✓ Numerically stable ($\delta = 0.01$)

Adaptive Cosine Schedule $\alpha(t)$

High guidance early (structure) \rightarrow fades to prompt conditioning.

- ✓ Smooth, differentiable
- ✓ Phase param ϕ controls scope



Gradient Clipping

$$\|d_{OT}\|_2 \leq \tau = 10.0$$

Prevents OT from overwhelming RF velocity.

Contributions

- (i) OT distances \leftrightarrow RF editing quality
- (ii) **Transport-Guided Inversion** with cosine scheduling
- (iii) **Transport-Enhanced FlowEdit**
- (iv) Validated on **FLUX & SD3**
- (v) **Training-free** — zero computational overhead

Key Numbers

- $\times 135$ LPIPS — RF-Inversion+OTC
- $+16.4$ dB PSNR — FlowEdit (FLUX)+OTC
- -11.2% face distance — semantic editing
- -12.9% L2 — LSUN-Church stroke

Take-Away

Optimal Transport provides a **unified, training-free, $O(n)$** framework for rectified flow image editing across *both paradigms*.

Code & Results

github.com/marianlupascu/OT-RF

Thank you! Questions?

Algorithm 1: Transport-Guided RF Inversion (simplified)

Require: Source x_0 , prompt p , model v_θ , β_0 , ϕ

- 1: $z_0 \leftarrow \text{Encode}(x_0)$; $z_{\text{tgt}} \leftarrow z_0$; $z_T \leftarrow \text{RFInvert}(z_0, v_\theta)$
- 2: **for** $t = T$ **down to** 0, **step** $-\Delta t$ **do**
- 3: $v_{\text{tar}} \leftarrow v_\theta(z_t, t, p)$; $v_{\text{ref}} \leftarrow v_\theta(z_t, t, \emptyset, z_0)$
- 4: $v_{\text{RF}} \leftarrow v_{\text{tar}} + \eta(v_{\text{ref}} - v_{\text{tar}})$
- 5: $d_{\text{OT}} \leftarrow (z_{\text{tgt}} - z_t) / \max(T - t, \delta)$
- 6: $\alpha \leftarrow \beta_0 \cdot \frac{1}{2}(1 + \cos(\min(\frac{T-t}{T\phi}, 1)\pi))$
- 7: $v_{\text{enh}} \leftarrow v_{\text{RF}} + \alpha \cdot \text{clip}(d_{\text{OT}}, \tau)$
- 8: $z_{t-\Delta t} \leftarrow z_t + \Delta t \cdot v_{\text{enh}}$
- 9: **end for**
- 10: **return** $\text{Decode}(z_0)$