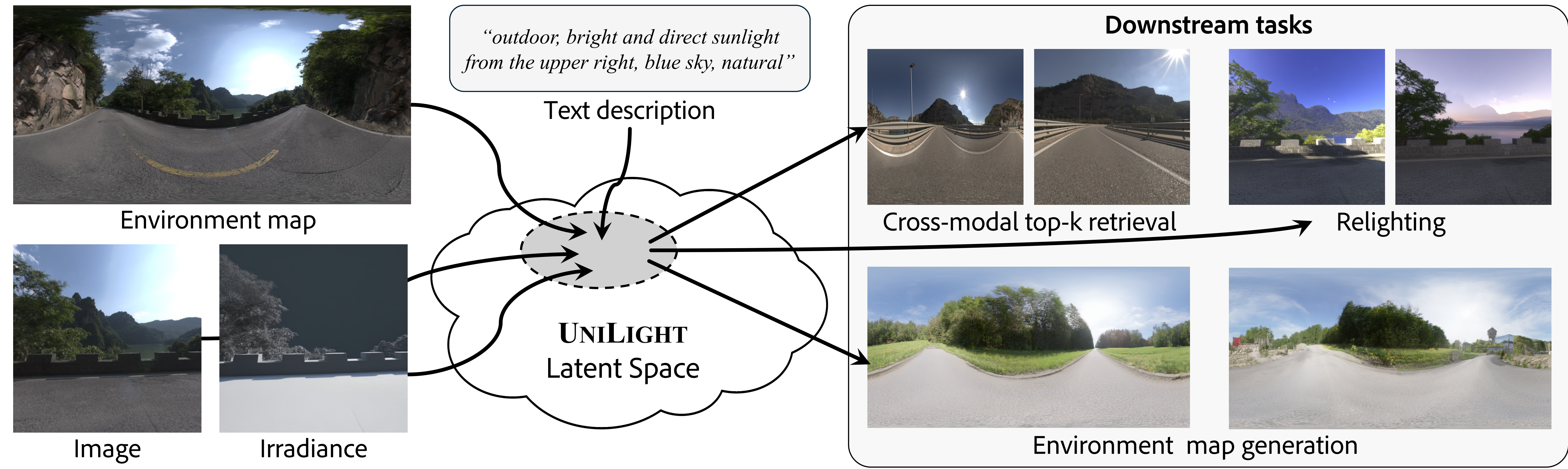


TL;DR

How can we represent lighting in images?

Encode **different light modalities** into a **joint latent space**, that can be used for retrieval, relighting and env. map generation.



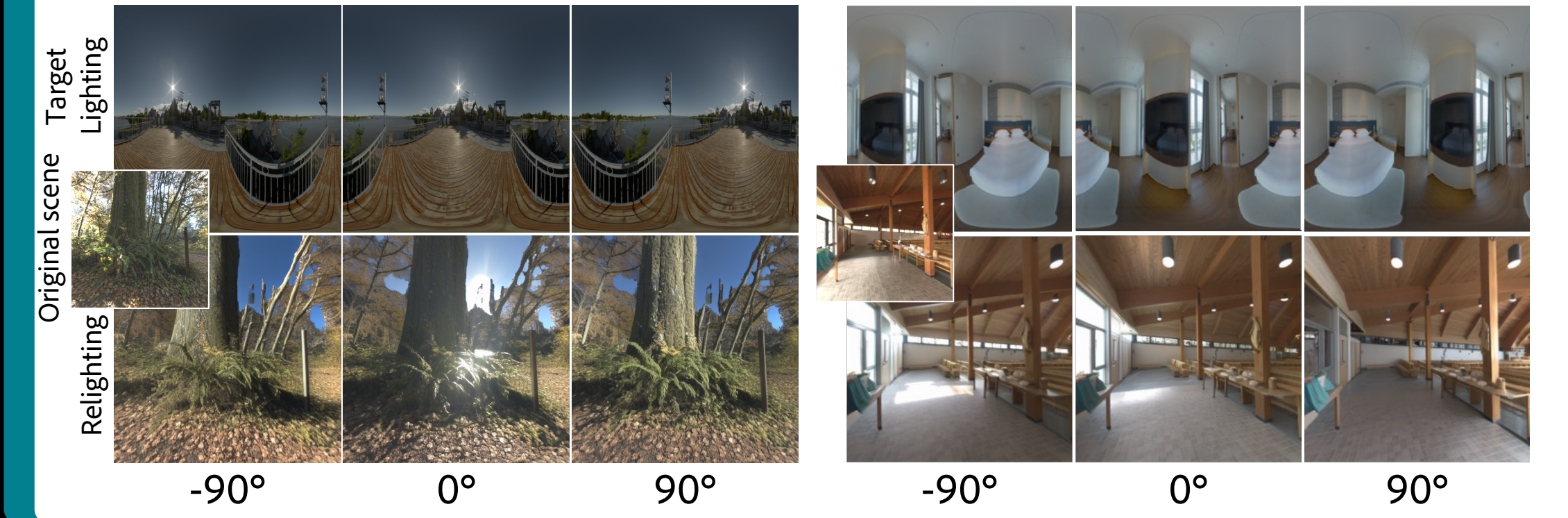
RETRIEVAL

Cross-modality retrieval based on the cosine similarity between latents.

Query	Top 1	Top 2	Top 3	Worst 1	Worst 2	Worst 3
Image → Env. map	0.996	0.996	0.960	-0.562	-0.535	-0.511
Text → Env. map	0.991	0.976	0.972	-0.584	-0.543	-0.541
Image → Text	0.998	0.997	0.993	-0.589	-0.513	-0.512

RELIGHTING

Fine-tune an X→RGB model to receive original scene intrinsics (depth, albedo and normals), and UniLight latents as conditioning for relighting.



ENVIRONMENT MAP GENERATION

Fine-tune a diffusion model to receive UniLight latents as conditioning, to generate env. map from either image, irradiance, or text descriptions.

Metrics on the rendering →

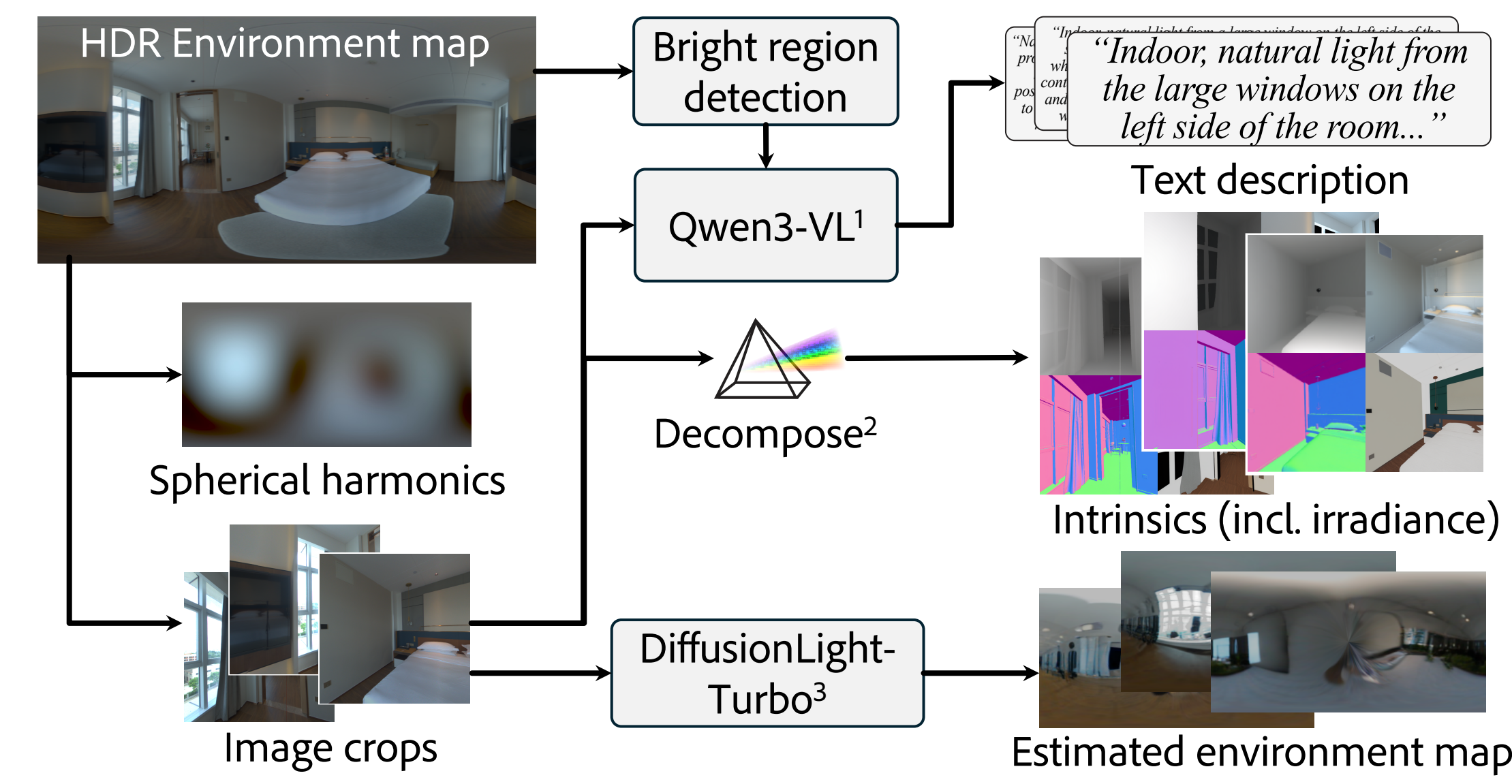
Method	PSNR _↑	RMSE _↓	SI-RMSE _↓	SSIM _↑	MAE _↓	LPIPS _↓
DiffusionLight-Turbo	27.77	0.157	0.062	0.902	0.148	0.088
UniLight	28.85	0.133	0.060	0.915	0.124	0.079

Visualization of the env. map and rendering ↓

The visualization shows a grid of environment maps and their corresponding renderings. The columns are labeled 'Input image', 'Ours', 'DiffusionLight-Turbo', and 'Ground truth'. Each row shows a different scene, with the 'Ours' column showing results that are more similar to the 'Ground truth' than the 'DiffusionLight-Turbo' results.

ALIGNED LIGHT MODALITIES

Starting from an HDR env. map., we extract aligned image crops, spherical harmonics, text description, intrinsics, and estimated env. map.



ARCHITECTURE

Modality-specific encoders for env. maps, images, irradiance, and text are trained contrastively to align their representations, with an auxiliary spherical-harmonics prediction task reinforcing directional understanding.

