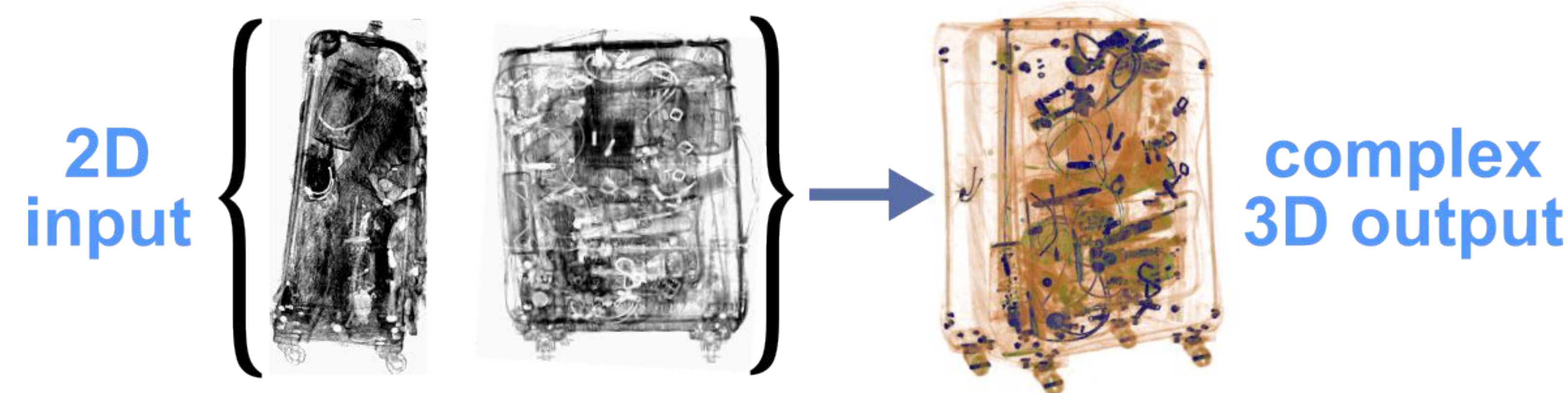


Unaligned 2D to 3D Translation with Conditional Vector-Quantized Code Diffusion using Transformers

Abril Corona-Figueroa, Sam Bond-Taylor, Neelanjan Bhowmik,
Yona Falinie A. Gaus, Toby P. Breckon, Hubert P. H. Shum, Chris G. Willcocks
Durham University, UK

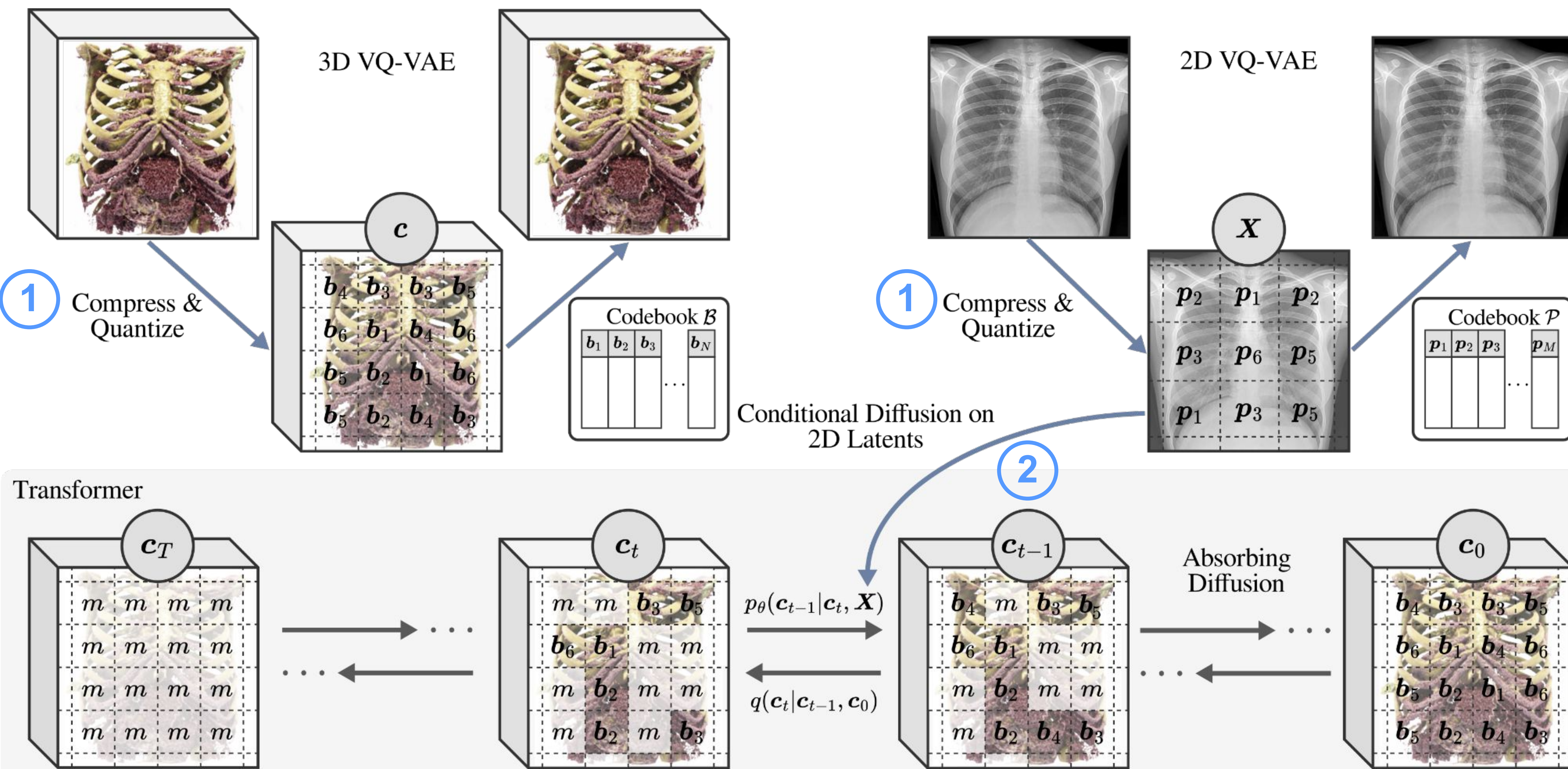
Motivation



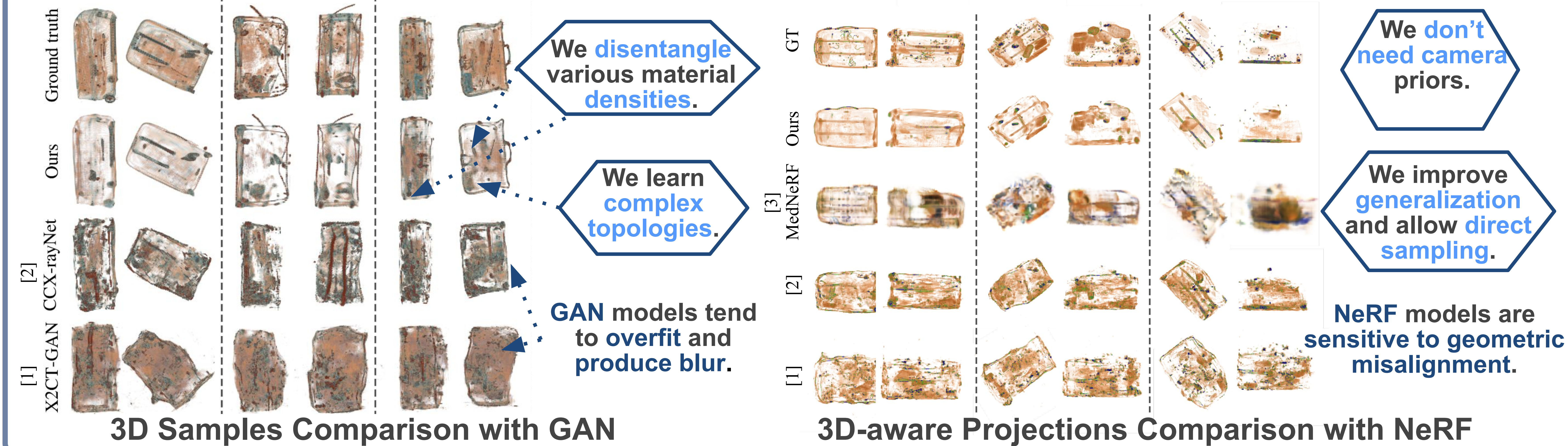
- Can we model **complex 3D** objects from two 2D images?
- **Unaligned inputs** for real-world applications.
- Improve **feature learning** and **speed up** 3D sampling.

Conditional Diffusion using Transformers

- Unpaired compression:** Learn rich Vector-Quantized discrete 2D and 3D spaces independently.
- Model $p(c|\mathbf{X})$ with a **diffusion model** parameterized with an **unconstrained transformer**.
 - \mathbf{c} : VQ codes of 3D data.
 - $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$: set of VQ codes of all 2D views.



Modeling Objects in Baggage Screening



Modeling Chest Structures



Quantitative Evaluation

Baggage Security Screening dataset

Method	↓ NLL	↑ Density	↑ Coverage	↑ SSIM	↑ PSNR
[1] X2CT-GAN	N/A	0.95	0.80	0.655	34.68
[2] CCX-rayNet	N/A	1.28	0.89	0.886	35.45
Ours	0.007	2.01	0.97	0.899	39.45

LIDC-IDRI (chest) dataset

Method	↓ NLL	↑ Density	↑ Coverage	↑ SSIM	↑ PSNR
[1] X2CT-GAN	N/A	0.87	0.88	0.321	19.68
[2] CCX-rayNet	N/A	1.41	0.98	0.386	22.66
Ours	0.10	1.42	0.97	0.436	25.05

Contributions

- ★ Full-coverage attention allows unaligned inputs.
- ★ Domain-invariance (e.g. imaging modalities).
- ★ Fast high-quality 3D samples.
- ★ Data likelihood representation.
- ★ Generative control (e.g. feature level).
- ★ Mode coverage (i.e. 2D & 3D distributions).
- ★ Global context of conditional 2D inputs.

Not requisite of:

- ❑ Camera priors like in NeRF.
- ❑ Local alignment of inputs like in CNNs.
- ❑ Many input views (i.e. 2 views suffice).
- ❑ Continuous latent representations.
- ❑ Hierarchical architectures.
- ❑ Deep architectures with skip-connections.

