# Thinking Outside the Pool:
# Active Training Image Creation for Relative Attributes

Aron Yu[†]

Kristen Grauman[†‡]

[†]*University of Texas at Austin*

[‡]*Facebook AI Research*

facebook research

## Fine-Grained Visual Comparisons

**More? Less? Equal?**

o compares images exhibiting *subtle visual differences* w.r.t. a target attribute

o trains attribute-specific ranking models using pairwise labels
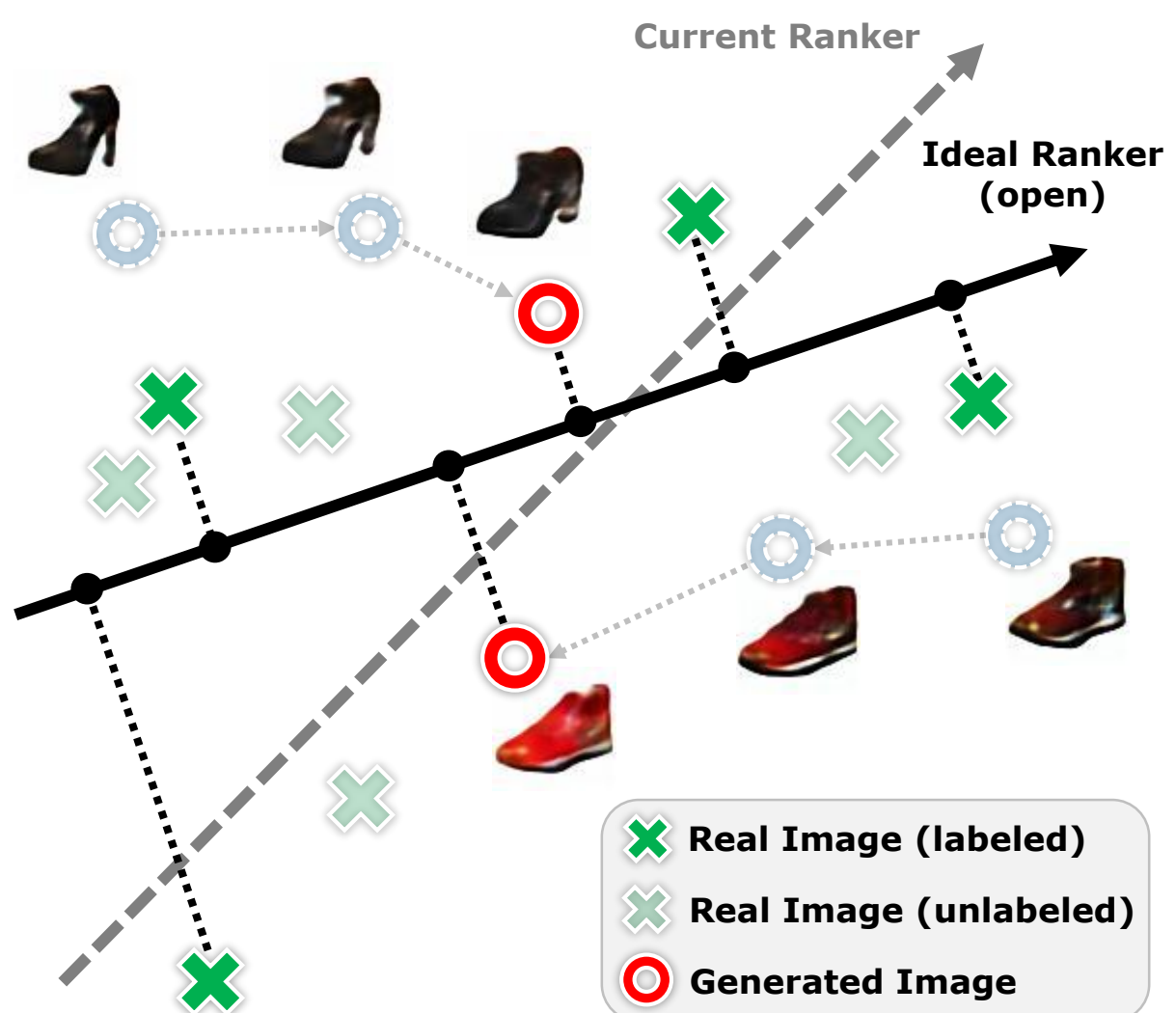
Sporty | vs.
Natural | vs.

**Curation Limit**

o visual variety and information from web content reach an upper limit

o subtle differences hard to directly curate in large numbers

**Pool-Based Active Learning**

o suffers from the "streetlight effect"

o starts with an existing pool of images

[Freytag et al. '14, Vijayanarasimhan et al. '14, …]

THIS IS WHERE YOU LOST YOUR WALLET?

NO, I LOST IT IN THE PARK. BUT THIS IS WHERE THE LIGHT IS.

*Source: Wikipedia*

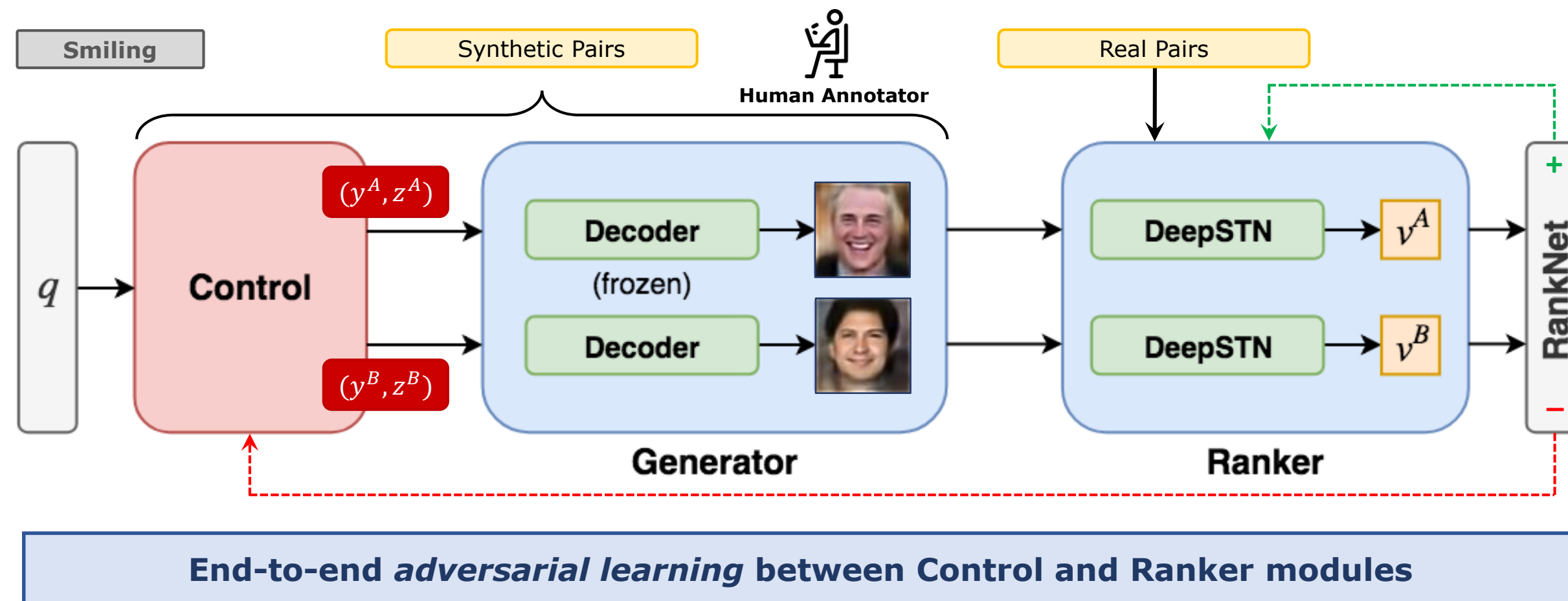## Think "Outside the Pool"



Current Ranker

Ideal Ranker (open)

**Our Idea**

o generate the most ideal training image pairs *directly*

o adversarial learning that allows the model to *actively teach itself*

o have human annotators label the generated image pairs

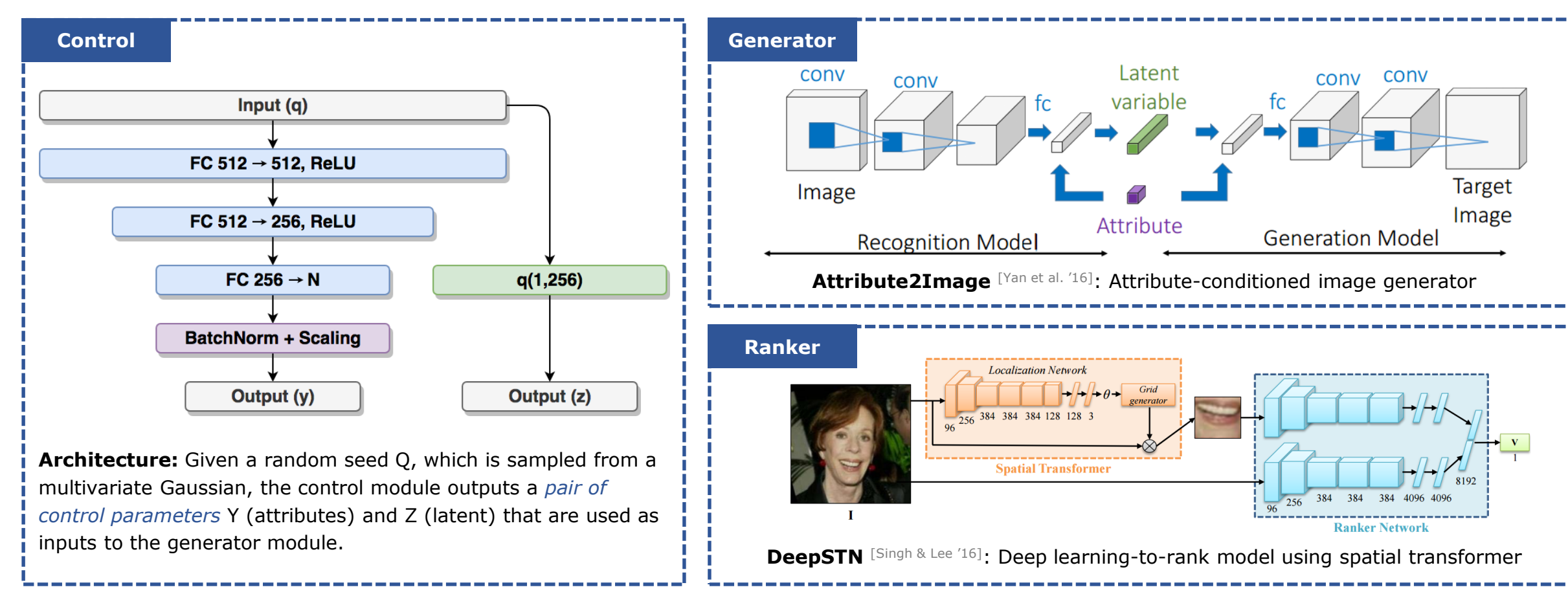o iteratively add to the pool of existing training image pairs

✖ Real Image (labeled)
✖ Real Image (unlabeled)
⬤ Generated Image

**Key Idea:** Joint learning of the visual task and training image generation
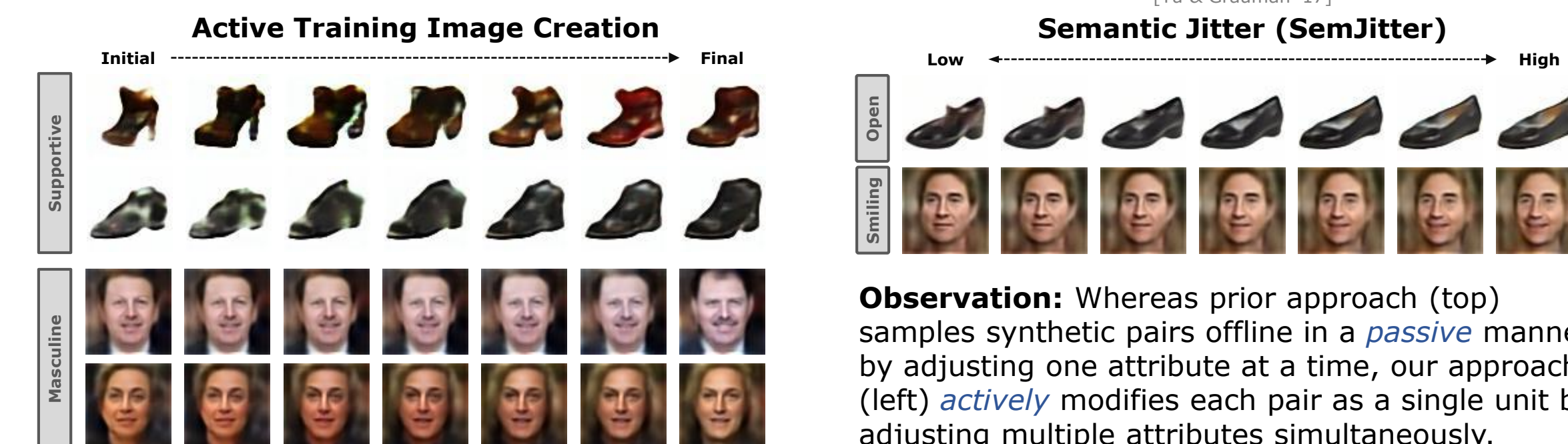
## AcTive Training Image Creation (ATTIC)



Smiling | Synthetic Pairs | Human Annotator | Real Pairs

$q$ → Control → $(y^A, z^A)$ → Decoder (frozen) → DeepSTN → $v^A$ → RankNet +

$(y^B, z^B)$ → Decoder → DeepSTN → $v^B$ → −

**Generator** | **Ranker**

**End-to-end *adversarial learning* between Control and Ranker modules**

**Control**

Input (q)

FC 512 → 512, ReLU

FC 512 → 256, ReLU

FC 256 → N | q(1,256)

BatchNorm + Scaling

Output (y) | Output (z)

**Architecture:** Given a random seed Q, which is sampled from a multivariate Gaussian, the control module outputs a *pair of control parameters* Y (attributes) and Z (latent) that are used as inputs to the generator module.
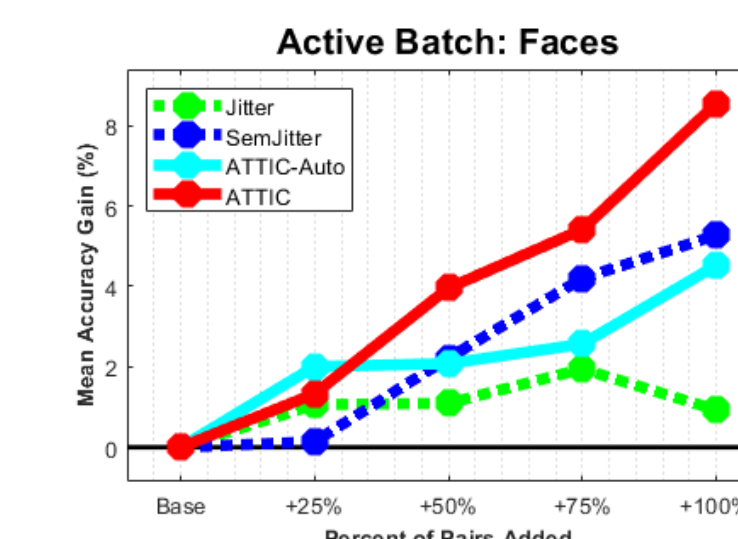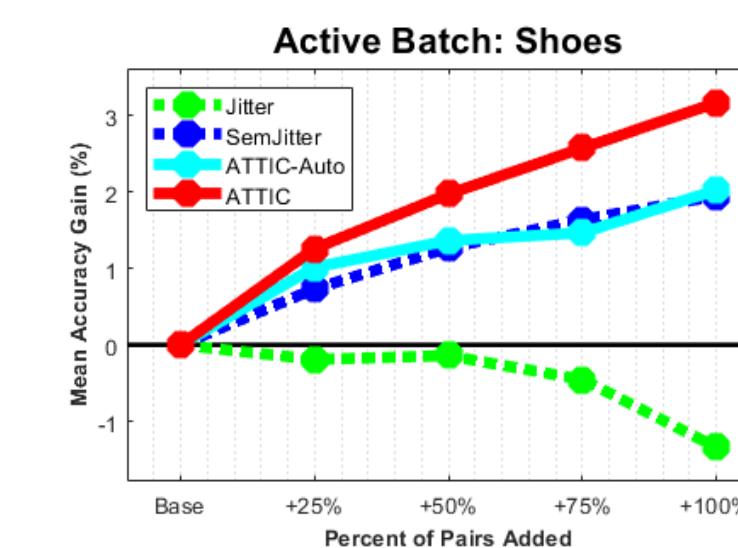
**Generator**

conv conv fc | Latent variable | fc conv conv

Image | Attribute | Target Image

Recognition Model | Generation Model

**Attribute2Image** [Yan et al. '16]: Attribute-conditioned image generator

**Ranker**

Localization Network | Grid generator | Ranker Network

Spatial Transformer

**DeepSTN** [Singh & Lee '16]: Deep learning-to-rank model using spatial transformer

## Image Synthesis & Progression



**Active Training Image Creation**
Initial ————— Final

Supportive
Masculine

[Yu & Grauman '17]

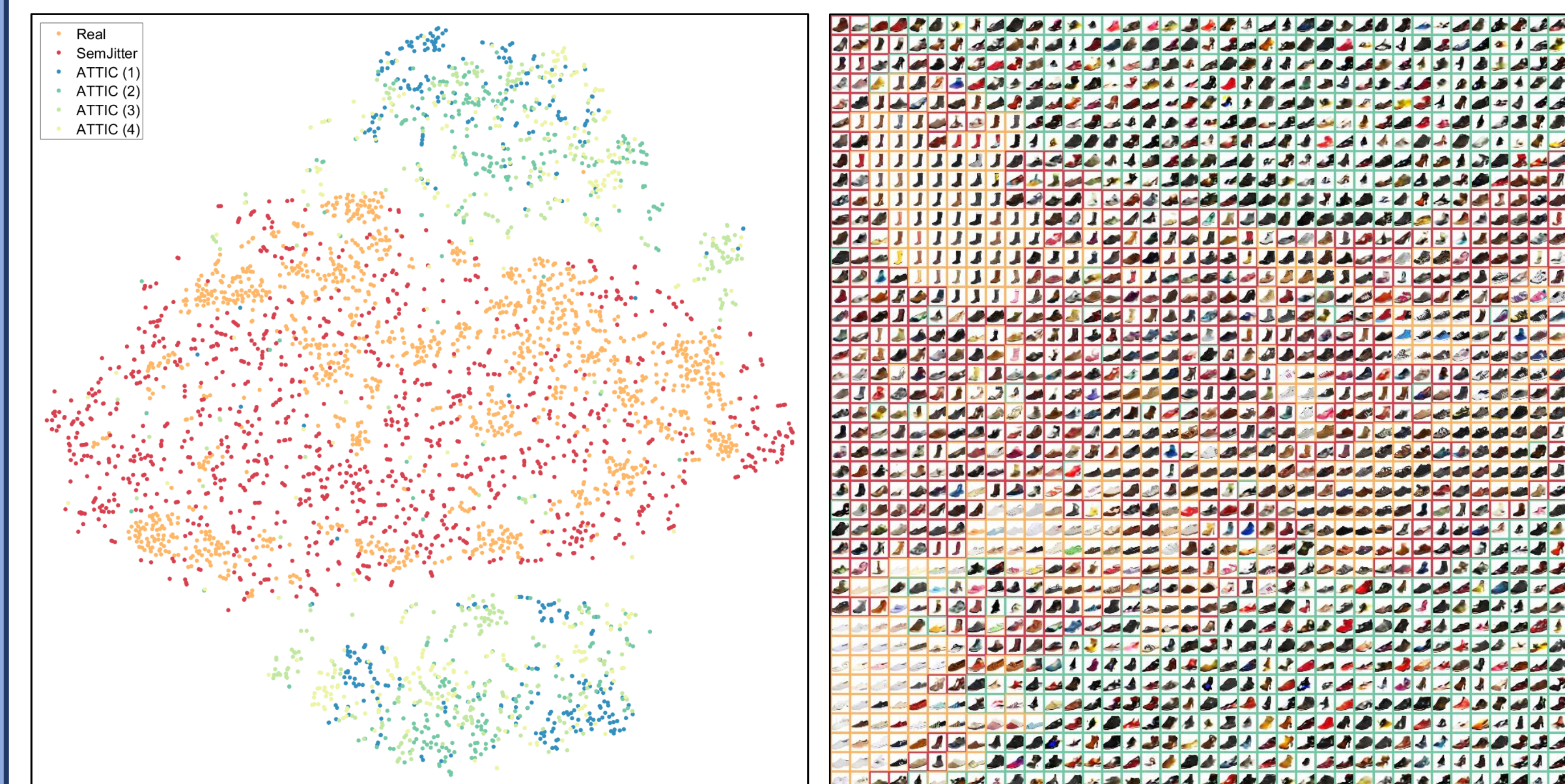**Semantic Jitter (SemJitter)**
Low ————— High

Open
Smiling

**Observation:** Whereas prior approach (top) samples synthetic pairs offline in a *passive* manner by adjusting one attribute at a time, our approach (left) *actively* modifies each pair as a single unit by adjusting multiple attributes simultaneously.

## Experimental Results



Real Pairs | Active Synthetic Pairs | Rejected

Casual | Durable | Smiling

Active Batch: Shoes
Active Batch: Faces

| | Real | Jitter | SemJitter | ATTIC | SemJitter (Auto) | ATTIC (Auto) |
|---|---|---|---|---|---|---|
| **Zappos50K (shoes)** | 86.74 | 85.60 | 86.87 | 87.62 | 87.59 | **89.07** |
| **LFW/LFW-10 (faces)** | 82.30 | 82.99 | 84.81 | **84.99** | 83.73 | 84.59 |



**t-SNE Embedding:** Densification of the training space through interpolation and extrapolation beyond the real training images.