

# Visual Anagrams and Applications

Ashwin Baluja, 4/18/24

# Visual Anagrams: Generating Multi-View Optical Illusions with Diffusion Models

Daniel Geng, Inbum Park, Andrew Owens

University of Michigan

Correspondence to: dgeng@umich.edu



a drawing  
of a penguin



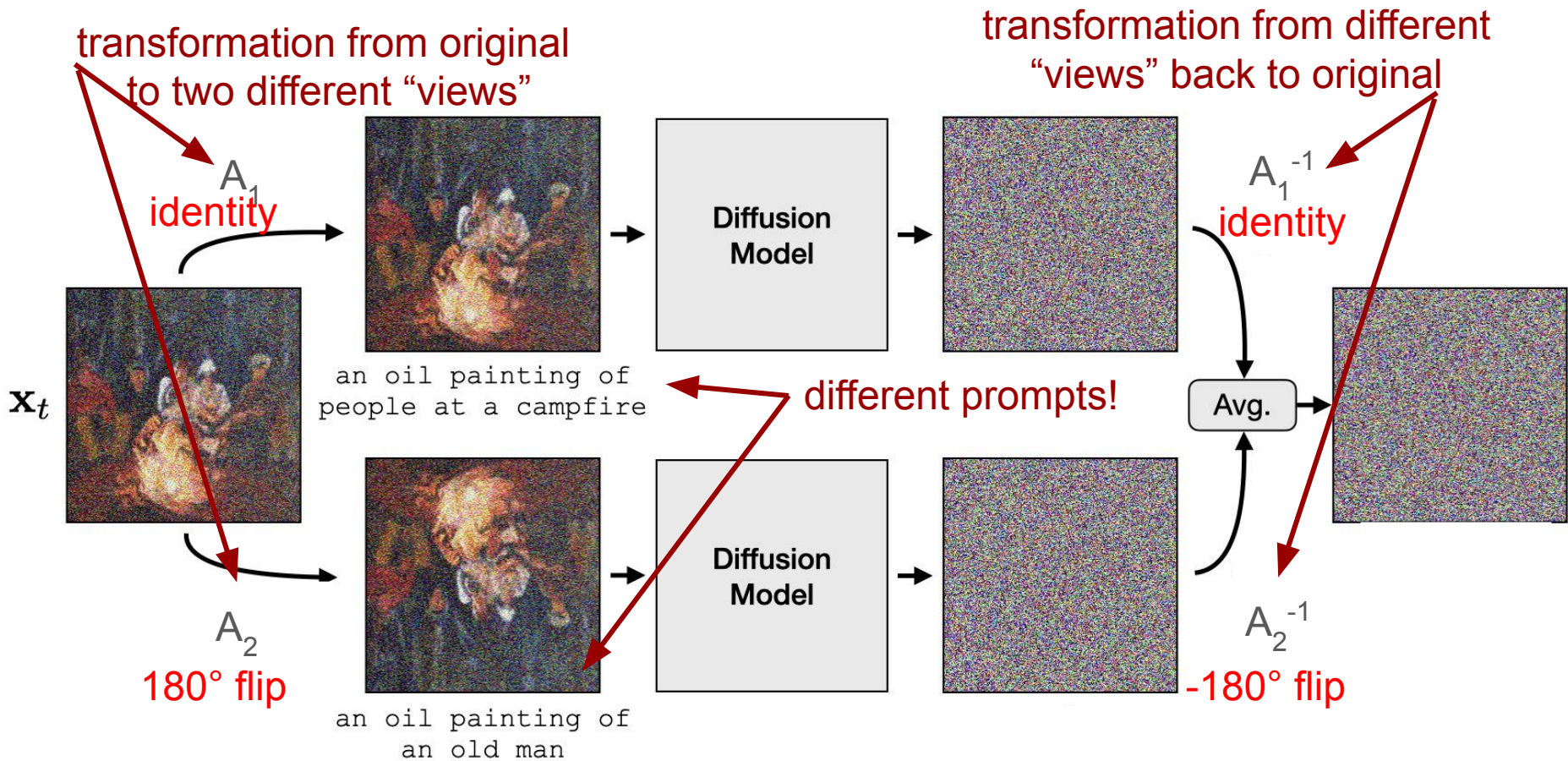
a drawing  
of a giraffe



an oil painting  
of a horse

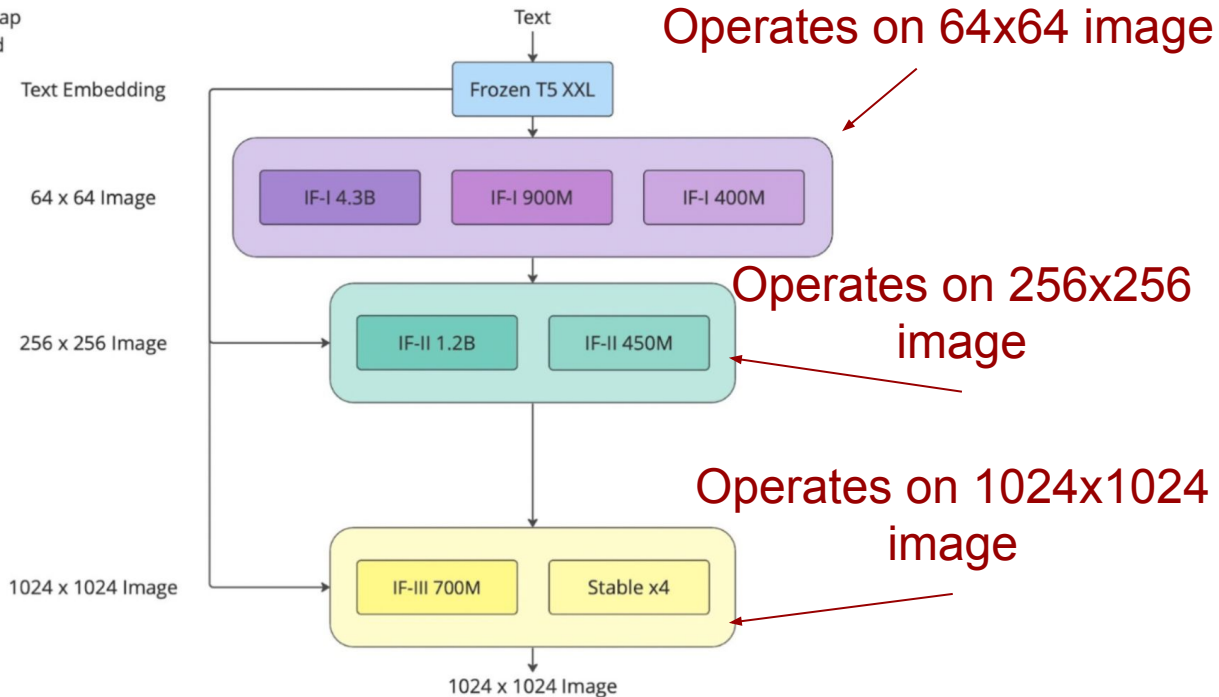


an oil painting of  
a snowy mountain village



# We need a direct relation between output pixels and noise

a photo of a violet baseball cap  
with yellow text "deep floyd  
better than text"



# What types of transformations work?

Expected input for  
diffusion model



Signal



$a_t$



Noise

transformation between  
different "views"

$$A \left[ \text{Signal} + a_t(\text{Noise}) \right] \equiv A \left[ \text{Signal} \right] + A \left[ a_t \text{ Noise} \right]$$

transformation back  
to original view

} A must be linear!

# What types of transformations work?



Diffusion model tries  
to predict noise



}  $\mathcal{N}(0, 1)$

$A$  (our transformation) must preserve  $\mathcal{N}(0, 1)$

$$A(\text{pred\_noise}) \sim \mathcal{N}(0, 1) \implies \text{Cov}(A(\text{pred\_noise})) = I$$

(given that mean = 0)  $\text{Cov}(A(\text{pred\_noise})) = AA^T$

$$AA^T \implies A \text{ is orthogonal matrix}$$

# What types of transformations work?

A is a linear matrix that is orthogonal



“flips, rotations, skews, color inversions, and jigsaw rearrangements”

“any orthogonal transformation works as a view with our method”

**View: Vertical Flip**

a photo of  
a wedding dress

a photo of  
an old woman



**View: Vertical Flip**

an oil painting of  
albert einstein

an oil painting of  
elvis



**View: Vertical Flip**

an oil painting of  
a red panda

an oil painting of  
a teddy bear



**View: Vertical Flip**

an oil painting of  
a kitchen

an oil painting of  
a botanical garden



**View: Vertical Flip**

a painting of  
a museum

a painting of  
a camel







Wedding dress

Old woman



Albert Einstein



Elvis Presley

**View: 90° Rotation**

a lithograph of a village in the mountains

a lithograph of a ship



**View: 90° Rotation**

a lithograph of a theater

a lithograph of a ship



**View: Negate**

an ink drawing of a red panda

an ink drawing of elvis



**View: Negate**

an ink drawing of waterfalls

an ink drawing of wine and cheese



**View: Negate**

a lithograph of waterfalls

a lithograph of a table





Village in the  
mountains



A ship



Ink drawing of  
waterfalls



Ink drawing of wine  
and cheese

# Key Takeaways

- Diffusion is flexible!
  - Manipulating noise **can still result in coherent outputs**
- Conditioning isn't the only way to incorporate info
  - This problem frequently is approached by blending prompt embeddings
  - “Conjoined” diffusion processes are a **conceptually simpler** way to do this!
  - **Blend noise instead!!**

# What properties do we need to apply this to other domains?

- a transformation that preserves diffusion properties...
  - (in the image case, a linear and orthogonal transformation matrix)
- and an obvious place to separate the task into separate diffusion processes
  - (in the image case, diffusing two images and averaging the noise)

# We can apply this to graphs!

At each step, a graph neural network:

- Aggregates information from each node's neighborhood by averaging
  - (in a graph attention network, a weighted average of neighbors, given by softmax of attention scores)
- Transforms aggregated info with a neural network
- Replaces each node's state with the transformed, aggregated info



# We can apply this to graphs! (modified for diffusion)

At each step, a graph neural network:

- **Diffuse over each node's state** (each estimate has mean of 0)
- Aggregates noise estimates from each node's neighborhood by averaging
  - (in a graph attention network, a weighted average of neighbors, given by softmax of attention scores)
  - **divide softmax'd attention scores by L2 norm to preserve variance of 1**
- ~~Transforms aggregated info with a neural network~~
- Replaces each node's state with the transformed, aggregated info

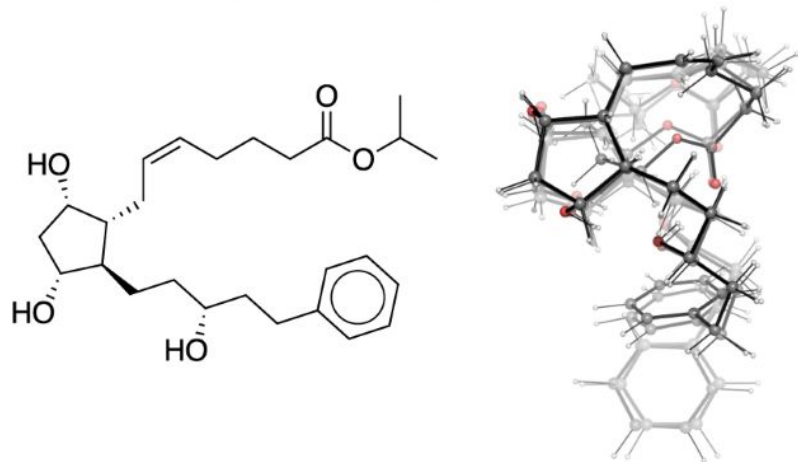
preserves diffusion properties... 

separate diffusion processes... 

# Problem Definition - Molecular Conformation

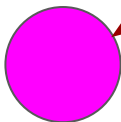
- Given a molecule graph, return a stable configuration of 3D coordinates for each atom
- Conditional diffusion: conditioned in molecule graph, diffuse coordinates

CC(C)OC(=O)CCC/C=C\C[C@H]1[C@@H](O)C[C@@H](O)[C@@H]1CC[C@@H](O)CCc1ccccc1



# What data do we use?

an atom



=

{ x, y, z, atom properties, molecule properties }

coordinates

per atom, per  
molecule

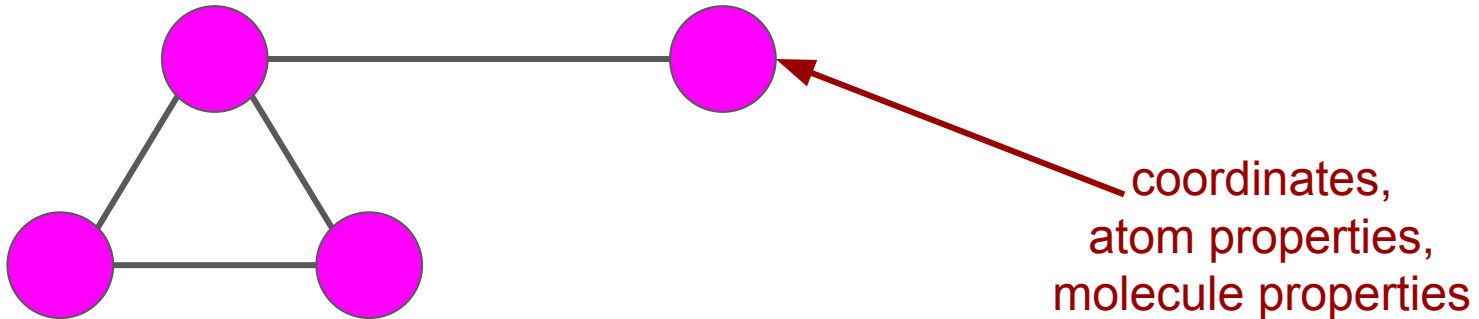
same for all atoms  
per molecule

from qm9

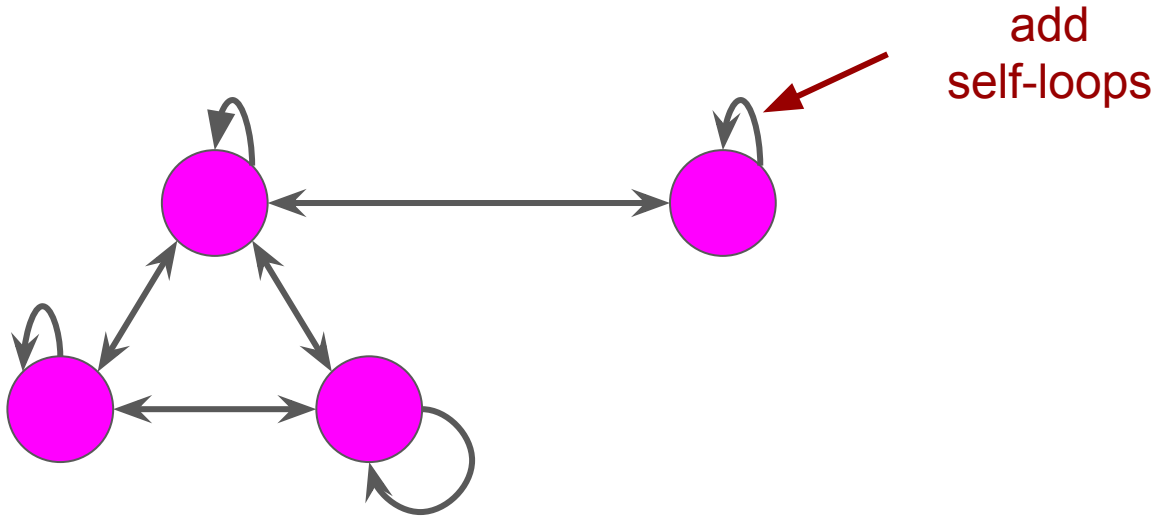
from MPNN for  
quantum chemistry

from qm9

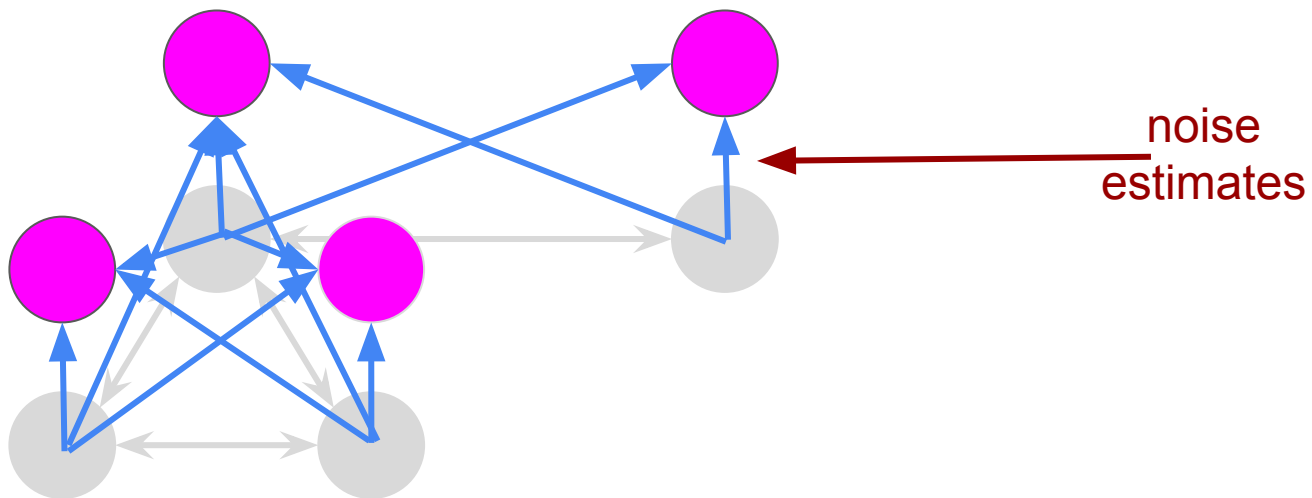
An example molecule:



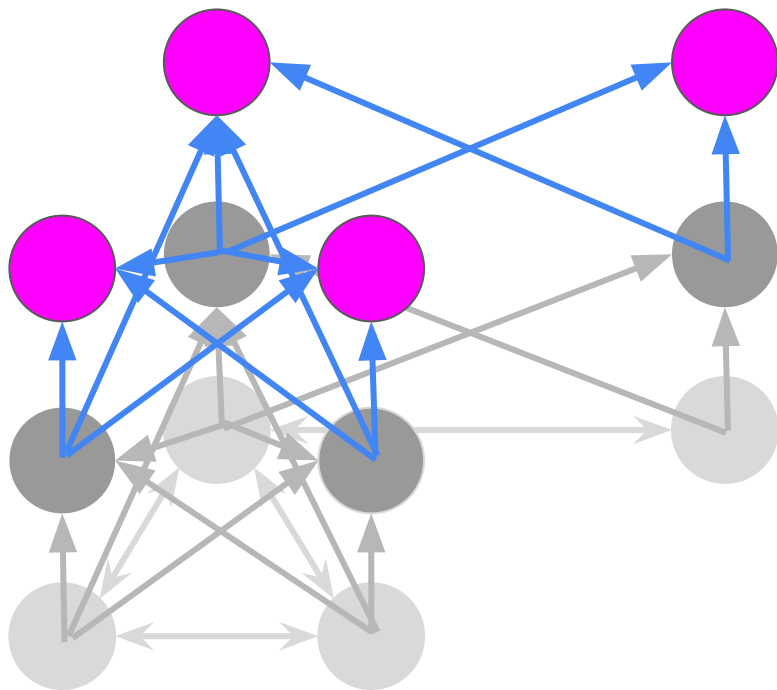
An example molecule: (at diffusion step 1)



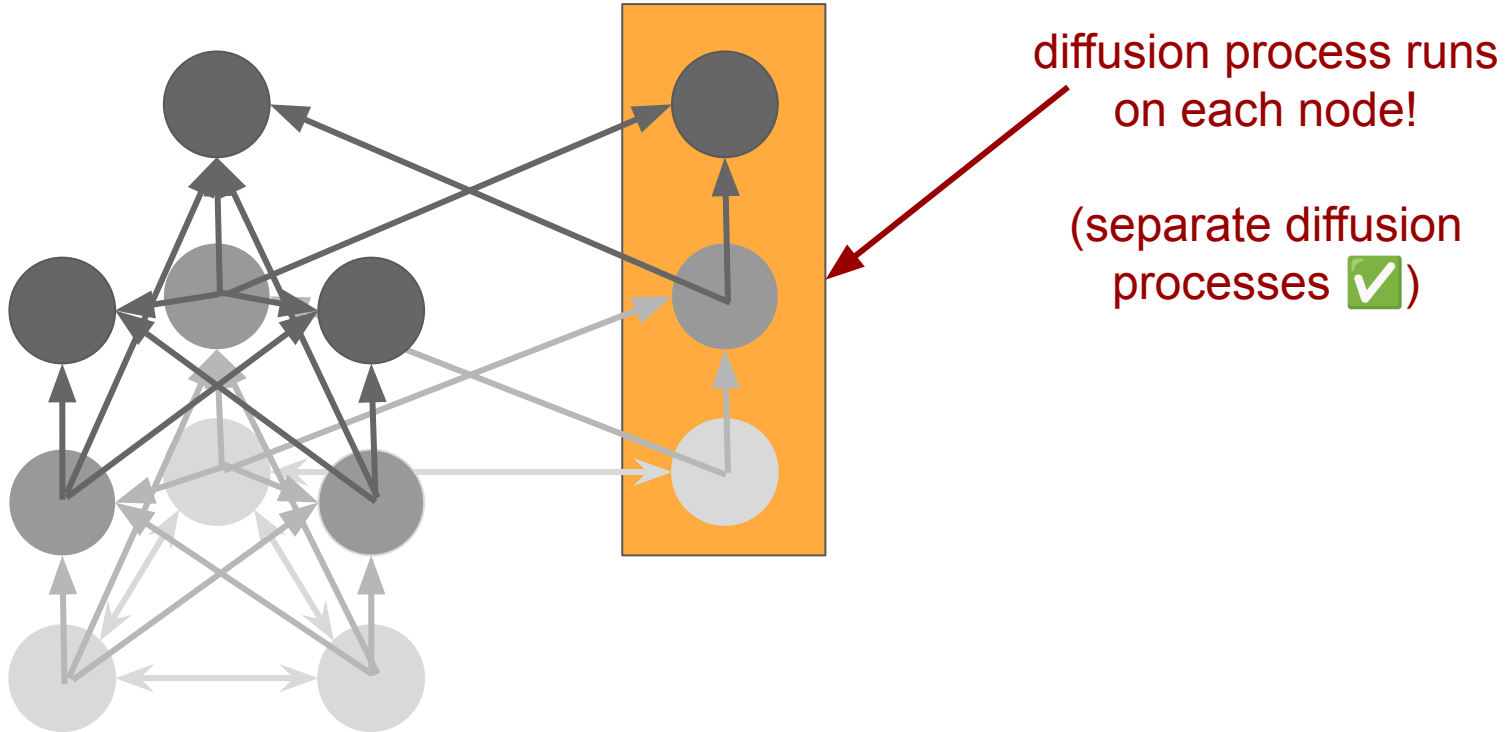
An example molecule: (at diffusion step 2)



An example molecule: (at diffusion step 3)

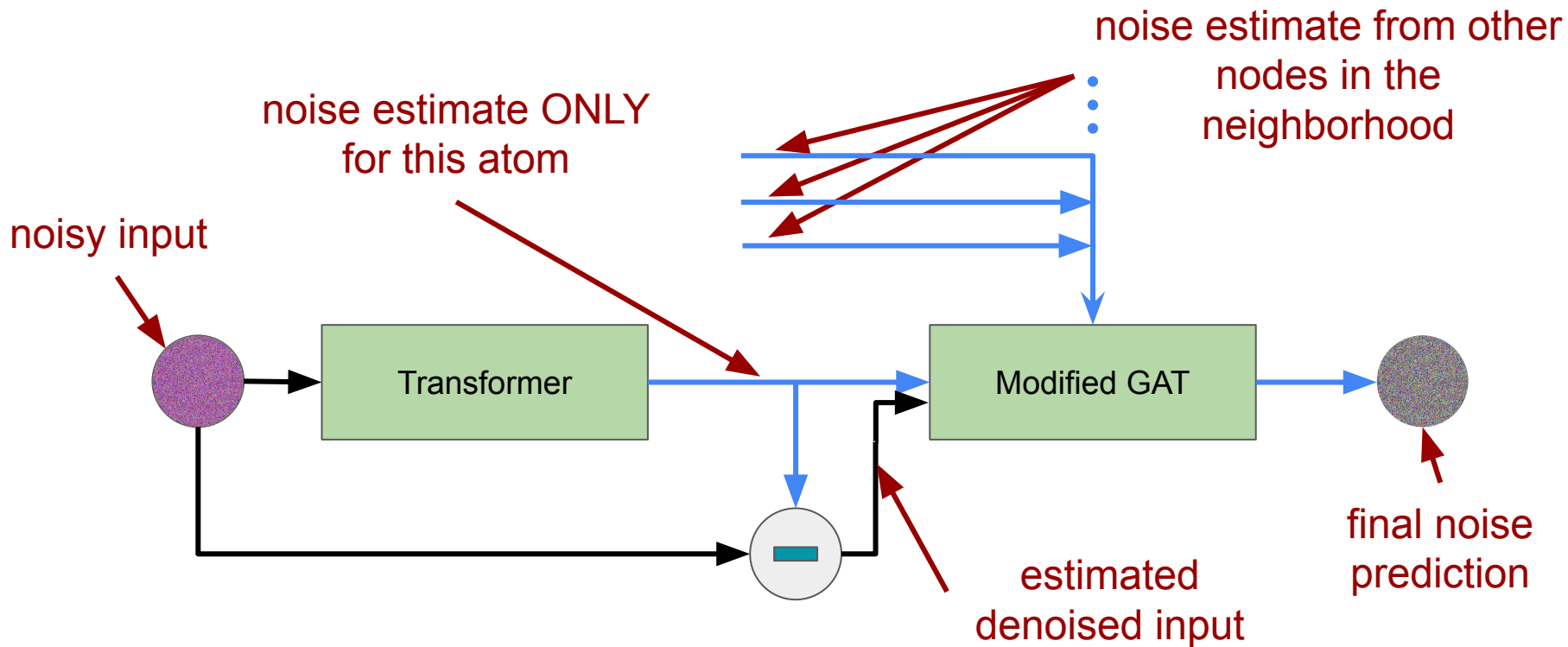


# An example molecule: explained

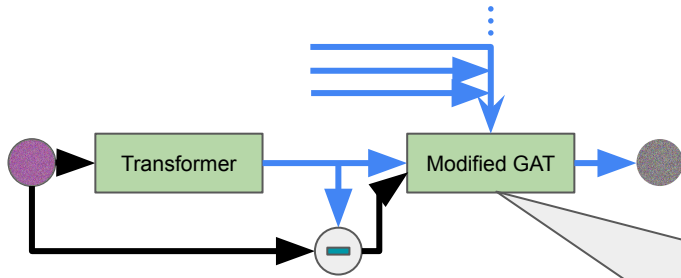




At each timestep, we...



# At each timestep, we...



preserve  
variance!

preserve  
mean!

- compute attention scores for each neighbor
- normalize attention scores s.t.  $\sum((att\_x)^2 \text{ for } x \text{ in neighborhood}) = 1$
- output =  $\sum(att\_x \text{ (noise\_est) for } x \dots)$

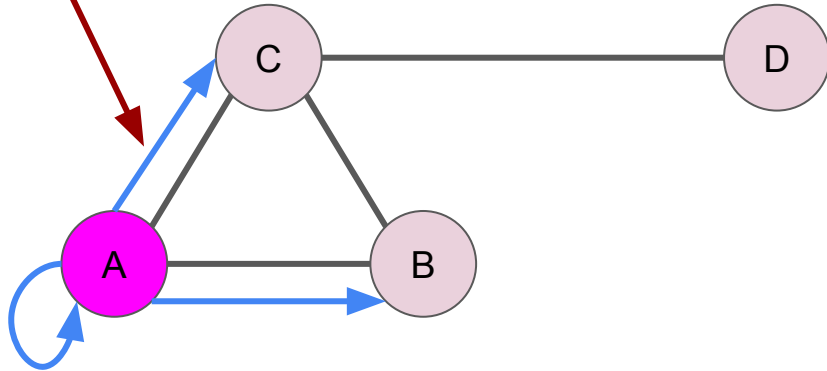
# Diffusion...

- A diffusion model diffuses through time
  - Each model step moves closer to the denoised signal
- A graph neural network propagates its information through time
  - Each GNN pass spreads information farther through the graph
- Each diffusion step only simulates a single timestep, but the method (diffusion) already takes that into account...

*How do we allow the information to propagate through time,  
while still only sampling individual timesteps?*

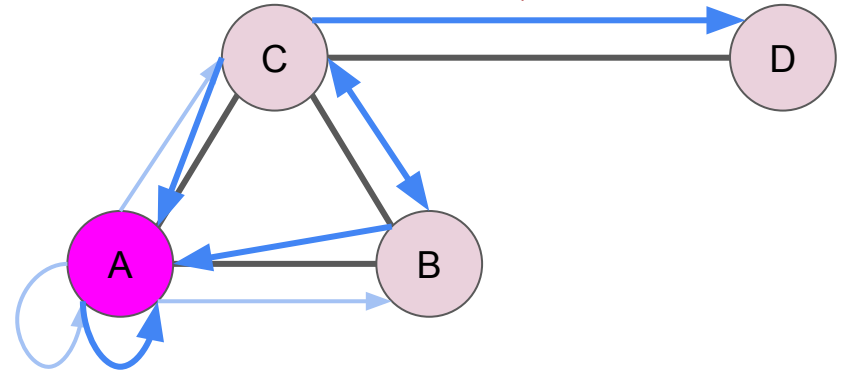
# Walks on a graph

Information ONLY passes along these edges



walks of length 1 starting at A

Looking only at one step, information from A will never make it to D



walks of length 2 starting at A

# Adjacency matrix trick...

- At timestep  $max\_timesteps - 1$ , information can only propagate to adjacent nodes
- At timestep 1, information will eventually propagate everywhere
- Lets simulate this propagation in one step:
  - At timestep  $t$  out of  $n$   $max\_timesteps$ , there are  $n - t$  timesteps left

All walks of length  $(n - t) = adjacency\_matrix^{(n-t)}$

this results in HUGE numbers, so we normalize adj a la GCN,  
 $(deg^{-1/2} adj deg^{-1/2})^{(n-t)}$

# All tied together now!

- train a diffusion process that fundamentally operates on individual atoms
- diffuse directly over the properties we are interested in
- blend the noise estimates together at each diffusion step
- preserve diffusion properties by normalizing attention coefficients and averaging 0-mean noise estimates
- simulate information propagation through time with matrix exponentiation

simpler, respects graph structure more

learns conformations by learning other useful properties

efficient!

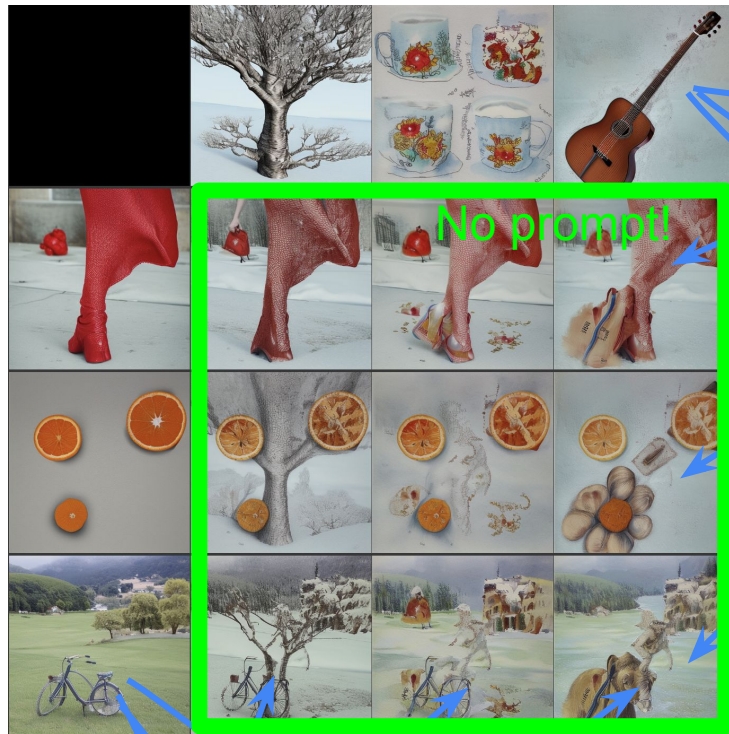
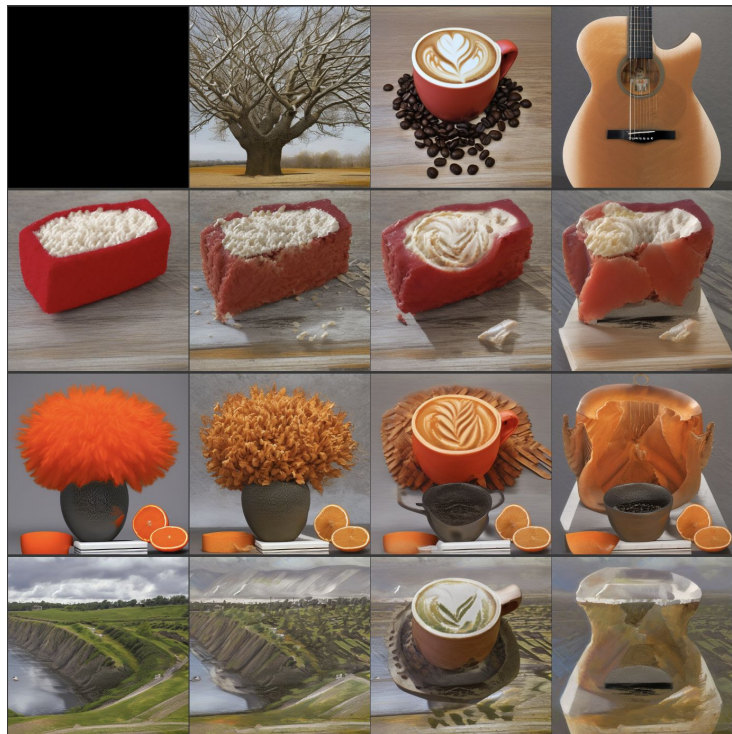
morally nice! blending noise estimates  
≅ every atom exerting a “force” on its neighbors, along bonds

simple correctness!

<https://github.com/ashwinbaluja/PerNodeDiffusion>

[https://colab.research.google.com/drive/1d\\_V2bsVZdtBwpOHHOt\\_WowH4U8Uu7GZn?usp=sharing](https://colab.research.google.com/drive/1d_V2bsVZdtBwpOHHOt_WowH4U8Uu7GZn?usp=sharing)

# Visual Anagrams ++ (latent-space) (my fun work) (extra)





# Visual Anagrams ++ (pixel-space) (my fun work) (extra)

