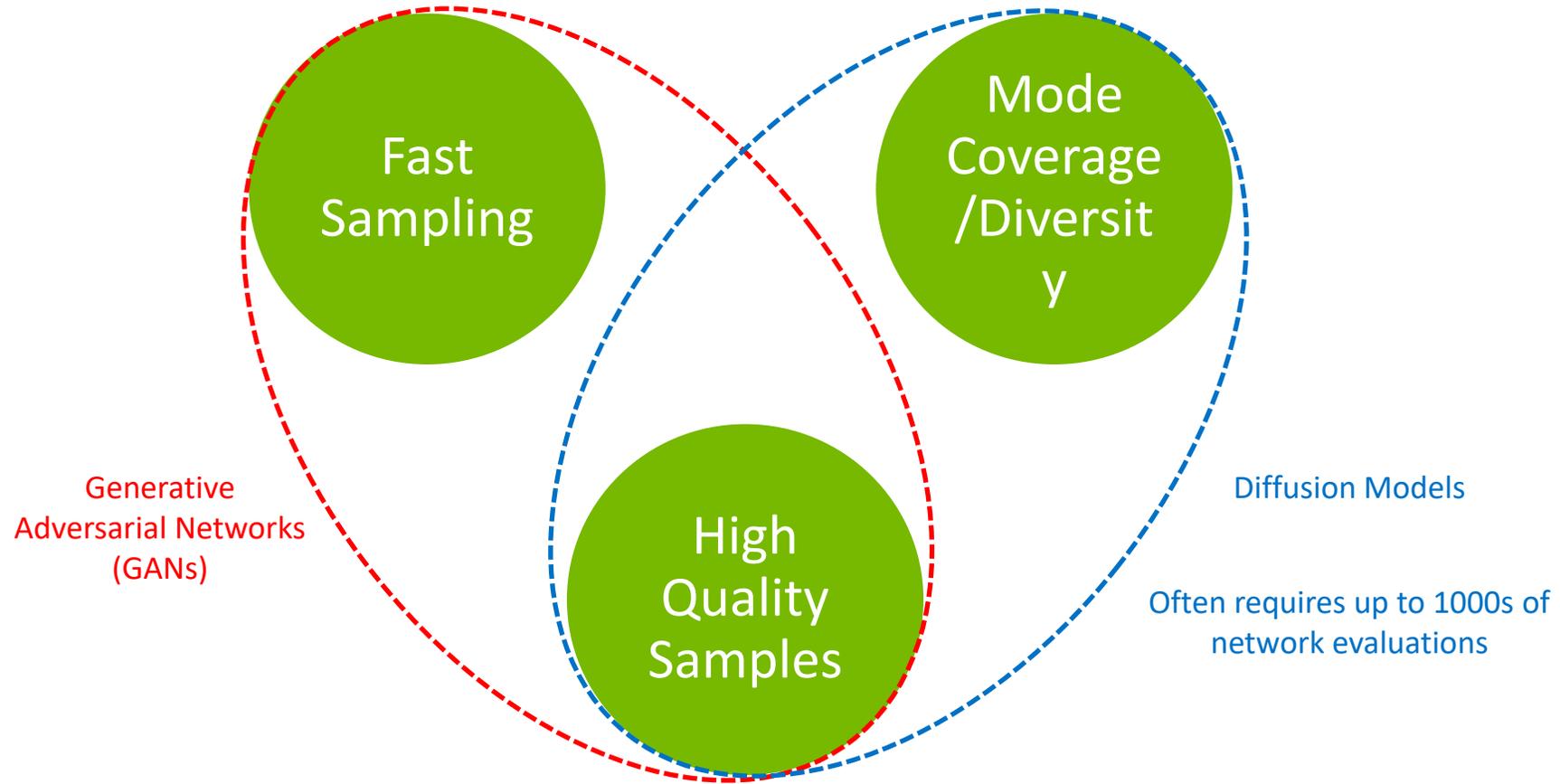


The background features a complex, abstract pattern of thin, overlapping lines in shades of green and white against a black background. The lines are arranged in a way that suggests motion and depth, with some lines forming a central, somewhat circular structure that resembles a stylized letter 'Q' or a similar shape. The overall effect is that of a high-speed, multi-layered visualization, possibly representing data flow or a neural network's internal state.

# Stable Diffusion Inference Optimization

Yiming Liu NVIDIA

# What makes a good generative model



# Sample Diversity

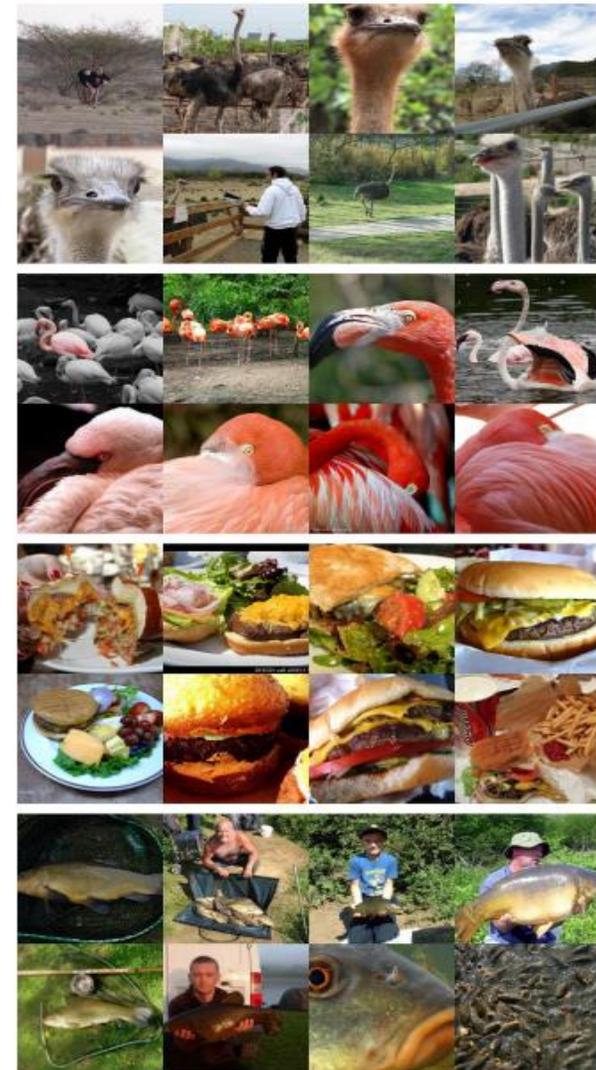
GANs



Diffusion

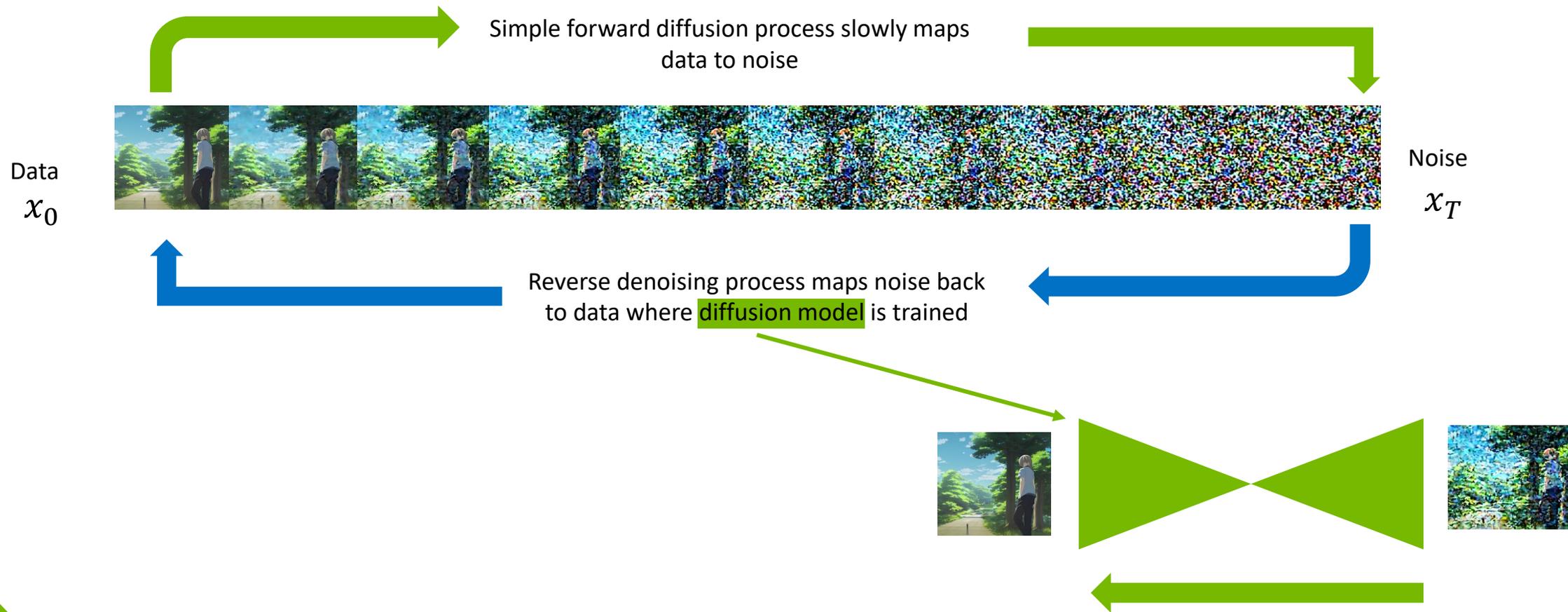


Real Data

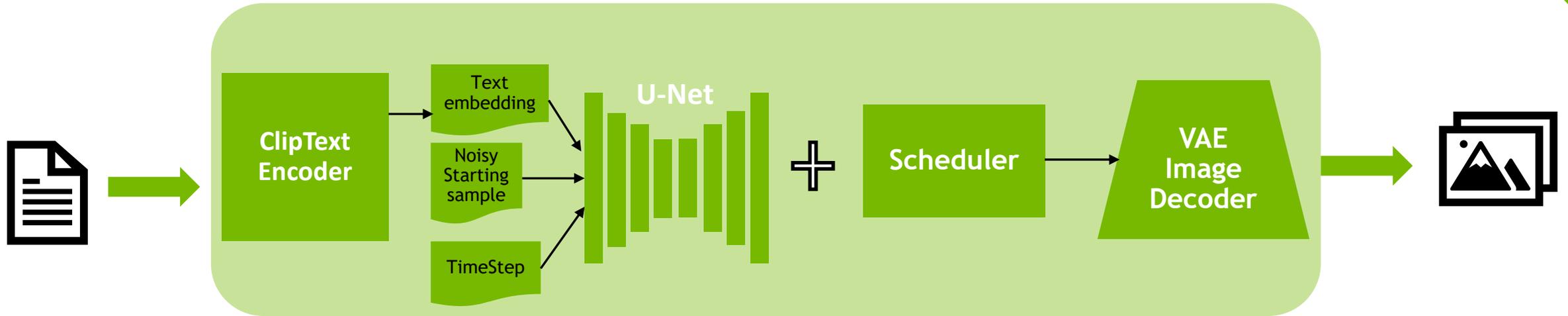


[Diffusion Models Beat GANs on Image Synthesis](#)

# Diffusion Process

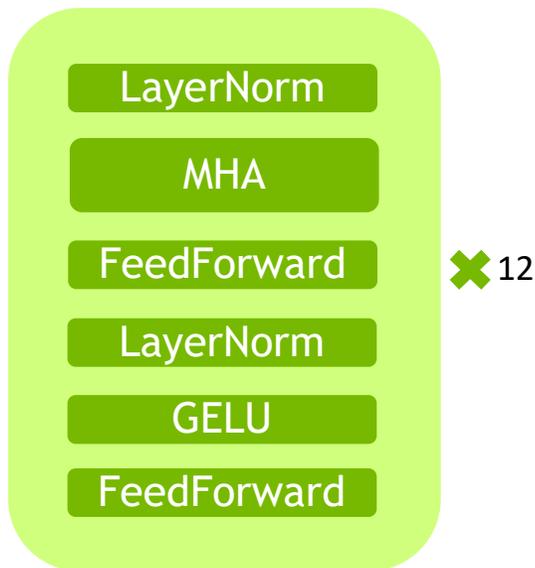


# Stable Diffusion Inference Pipeline



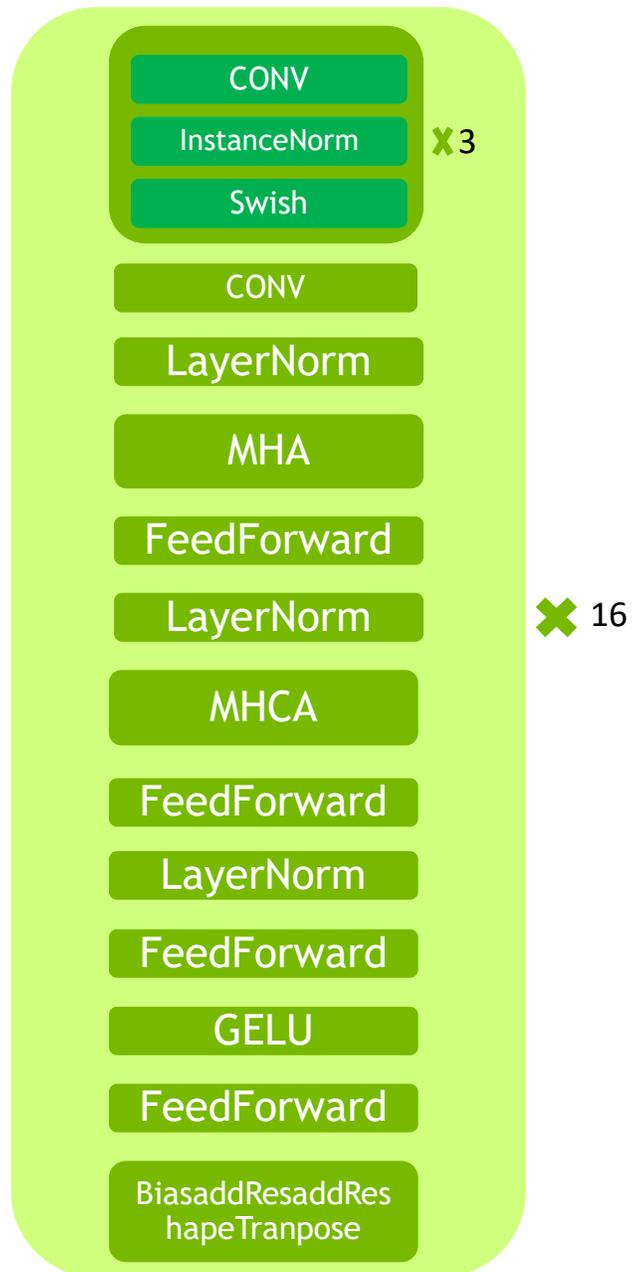
- **ClipText** for text encoding.  
Input: text.  
Output: 77 token embeddings vectors, each in 768 dimensions
- **U-Net + Scheduler** to gradually perform denoising in the information (latent) space.  
Input: text embeddings and a starting multi-dimensional array (structured lists of numbers, also called a *tensor*) made up of noise.  
Output: A processed information array
- **Autoencoder Decoder** that paints the final image using the processed information array.  
Input: The processed information array (dimensions: (4,64,64))  
Output: The resulting image (dimensions: (3, 512, 512) which are (red/green/blue, width, height))

## CLIP



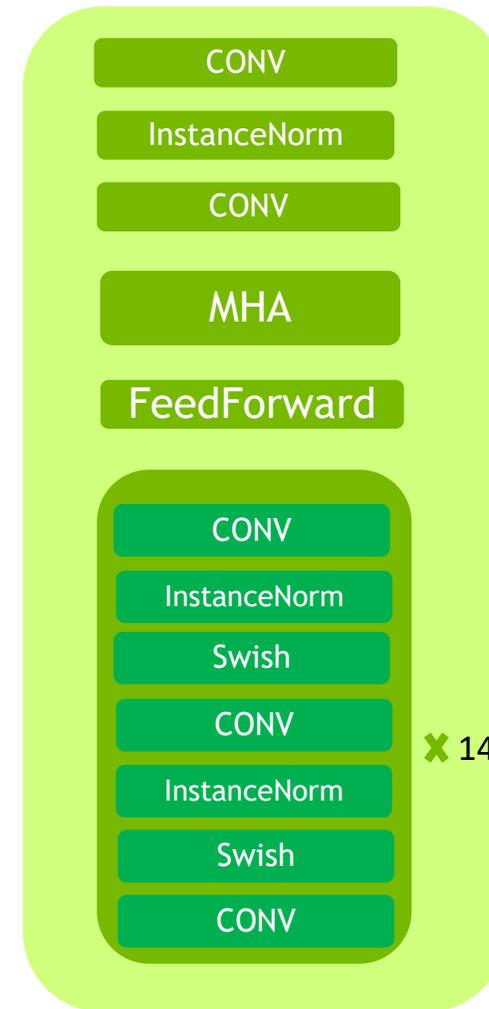
123M parameters

## U-NET



859M parameters

## VAE



49M parameters

Floating point operations measured under:

- $B=1$
- $H=W=64$

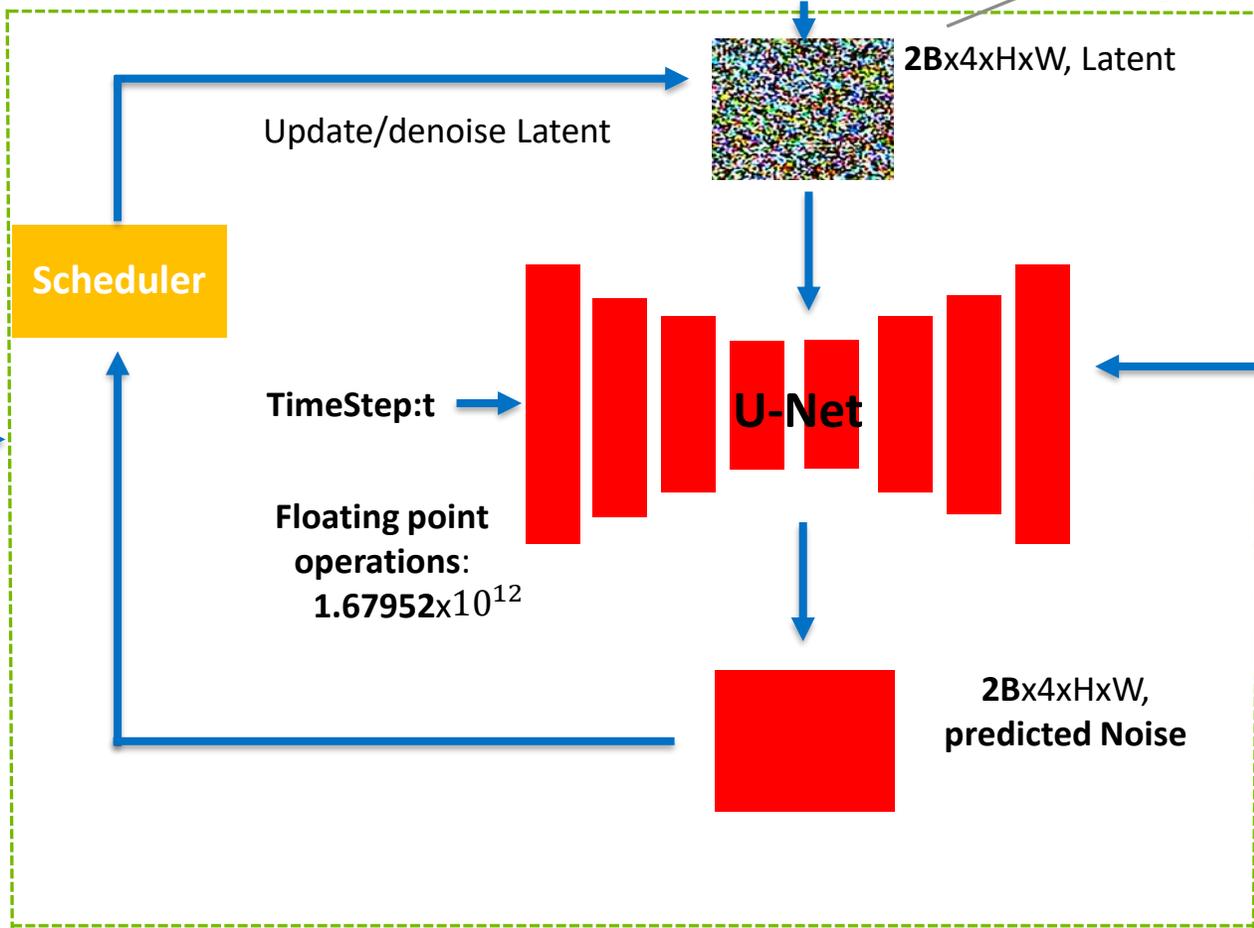
Starting Latent State: Sample Gaussian Noise from **random seed**

For classifier-free guidance

User prompt

Repeat  $T=50$  times

Floating point operations:  $83.976 \times 10^{12}$

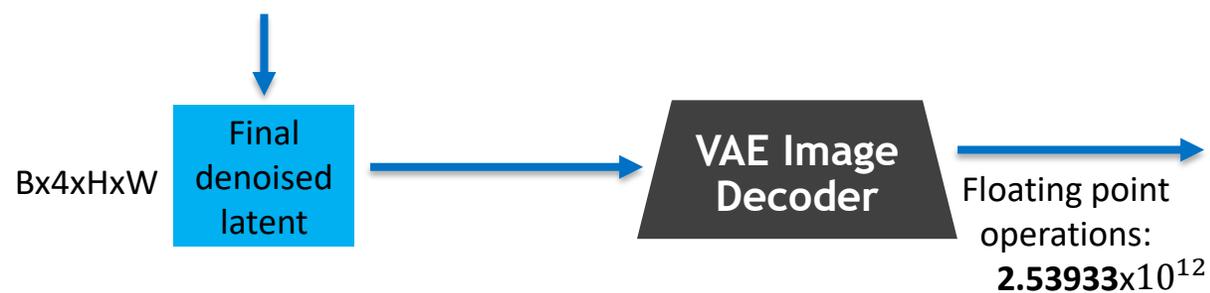


Floating point operations:  $0.02714 \times 10^{12}$

Clip text encoder

Negative prompt: Empty token if not provided

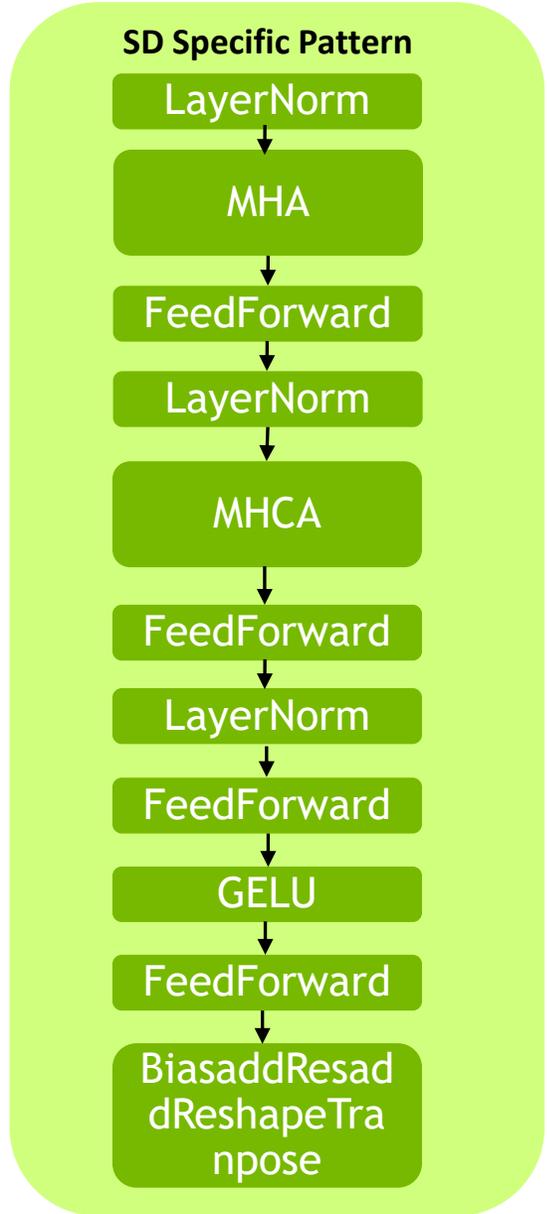
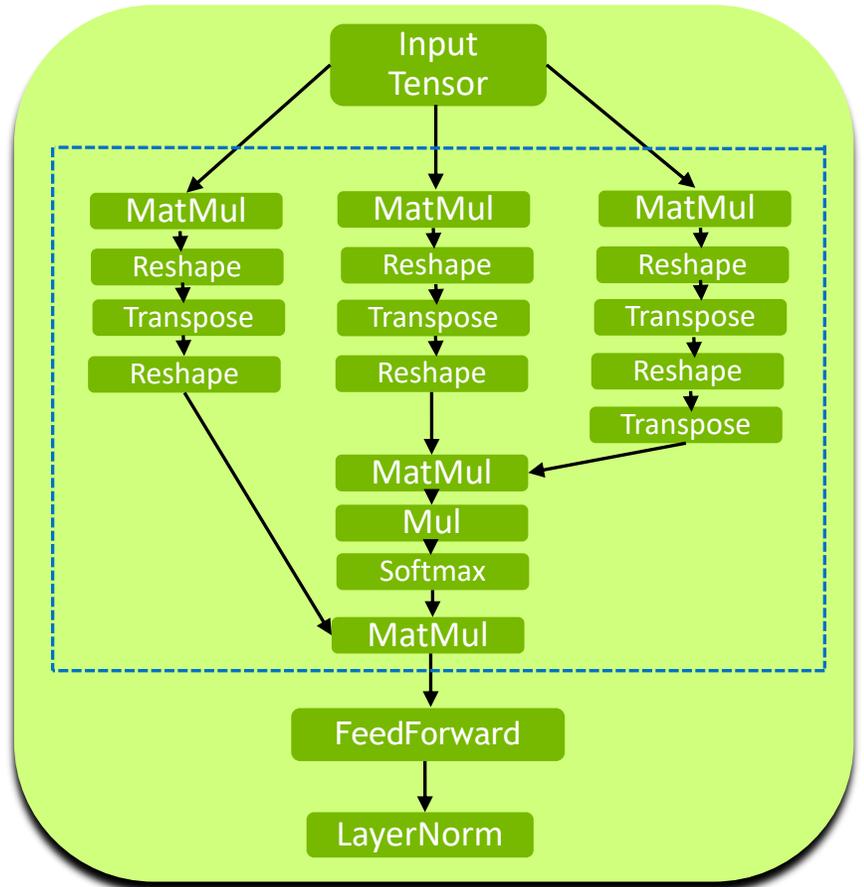
2Bx77x768, text embeddings



Output image: Bx3x8Hx8W



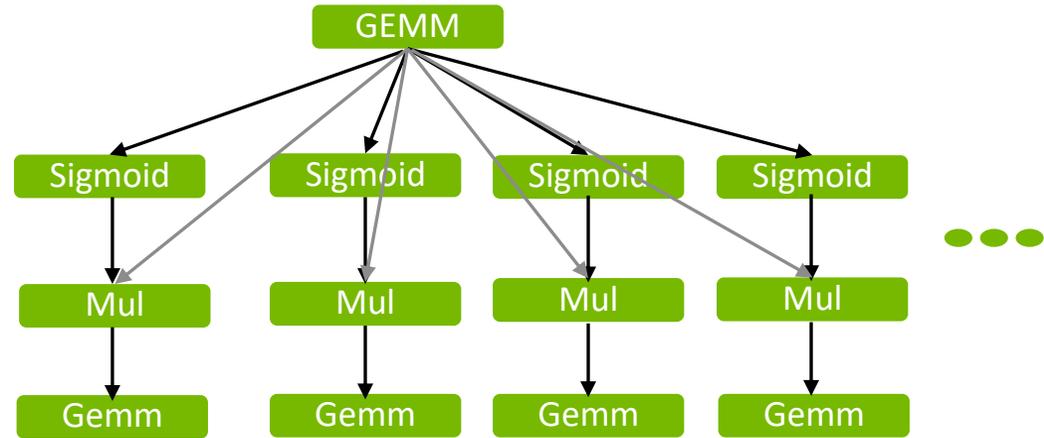
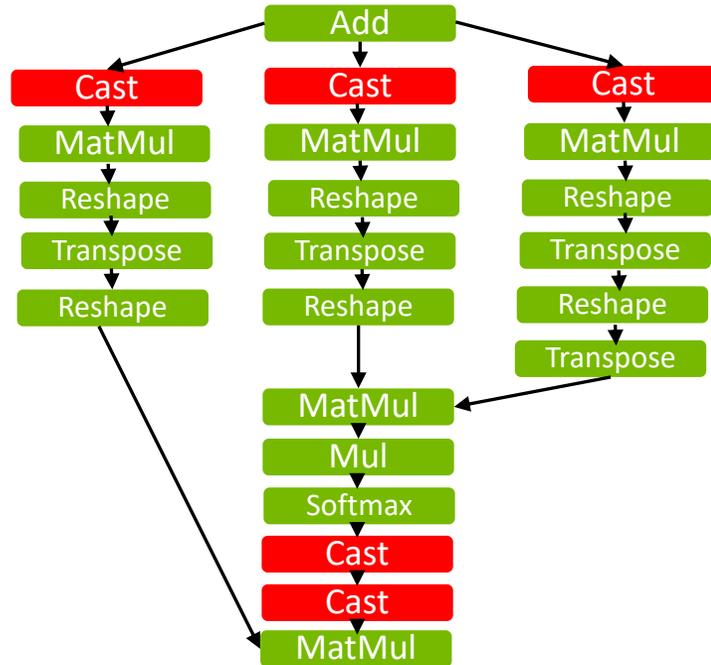
# How Does TensorRT Handle Attentions(since 8.0)



Myelin ForeignNode

- Detect the MHA pattern
- Try to offload MHA and its surrounding structures to **Myelin** as much as possible

# Utilize Myelin



- Remove redundant Cast onnx nodes to help TensorRT's **Myelin** backend compiler recognize the pattern as an MHA(not needed after TRT8.5.2)
- Merge repeated Swish ops

# Before Simplification

```
Dimensions": [1024,640,1,1],
"Format/Datatype": "Row major linear FP16 format"
}],
"TacticValue": "0x0000000000000000"
},{
"Name": "reshape_after_MatMul_5396 + Reshape_5419 + Transpose_5420",
"LayerType": "Shuffle",
"Inputs": [
{
"Name": "Reformatted Input Tensor 0 to reshape_after_MatMul_5396 + Reshape_5419 + Transpose_5420",
"Location": "Device",
"Dimensions": [1024,640,1,1],
"Format/Datatype": "Row major linear FP16 format"
}],
"Outputs": [
{
"Name": "onnx::Reshape_6685",
"Location": "Device",
"Dimensions": [1,8,1024,80],
"Format/Datatype": "Row major linear FP16 format"
}],
"ParameterType": "Shuffle",
"FirstTranspose": [0,1,2,3],
"Reshape": [1,1024,8,80],
"SecondTranspose": [0,2,1,3],
"ZeroIsPlaceholder": 1,
"TacticValue": "0x0000000000000000"
},{
"Name": "Reshape_5434",
"LayerType": "NoOp",
"Inputs": [
{
"Name": "onnx::Reshape_6685",
"Location": "Device",
"Dimensions": [1,8,1024,80],
"Format/Datatype": "Row major linear FP16 format"
}],
"Outputs": [
{
"Name": "query.171",
"Location": "Device",
"Dimensions": [8,1024,80],
"Format/Datatype": "Row major linear FP16 format"
}
```

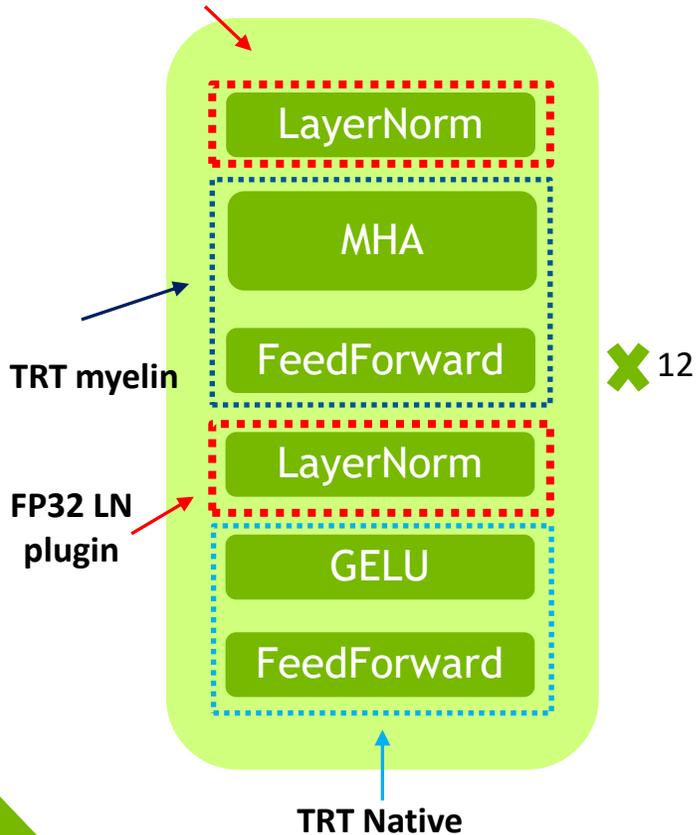
# After Simplification

```
"Activation": "NONE",
"HasBias": 1,
"HasReLU": 0,
"TacticName": "ampere_h16816cudnn_128x128_ldg8_relu_exp_stages_64x3_interior_nhwc_tn_v1",
"TacticValue": "0xdb058db40209ee1"
},{
"Name": "{ForeignNode[onnx::MatMul_9584 + (Unnamed Layer* 1263) [Shuffle]...Reshape_637 + Transpose_638]}",
"LayerType": "Myelin",
"Inputs": [
{
"Name": "onnx::Shape_808",
"Location": "Device",
"Dimensions": [1,320,64,64],
"Format/Datatype": "Channel major FP16 format where channel % 8 == 0"
}],
{
"Name": "encoder_hidden_states",
"Location": "Device",
"Dimensions": [1,77,768],
"Format/Datatype": "Row major linear FP16 format"
}],
"Outputs": [
{
"Name": "input.64",
"Location": "Device",
"Dimensions": [1,320,64,64],
"Format/Datatype": "Channel major FP16 format where channel % 8 == 0"
}],
"TacticValue": "0x0000000000000000"
},{
"Name": "InstanceBiasV-3",
"LayerType": "Constant",
"Inputs": [],
"Outputs": [
{
"Name": "(Unnamed Layer* 1564) [Constant]_output",
"Location": "Device",
"Dimensions": [1,32,1],
"Format/Datatype": "Row major linear FP16 format"
}],
"ParameterType": "Constant",
"weights": {"Type": "Float", "Count": 32,
"dimensions": [1,32,1],
```

To obtain the layer info json, you can build the engine with trtexec(<https://github.com/NVIDIA/TensorRT/tree/main/samples/trtexec>), with `--exportLayerInfo=layerinfo.json --profilingVerbosity=detailed` flags enabled.

# CLIP

FP32 LN plugin



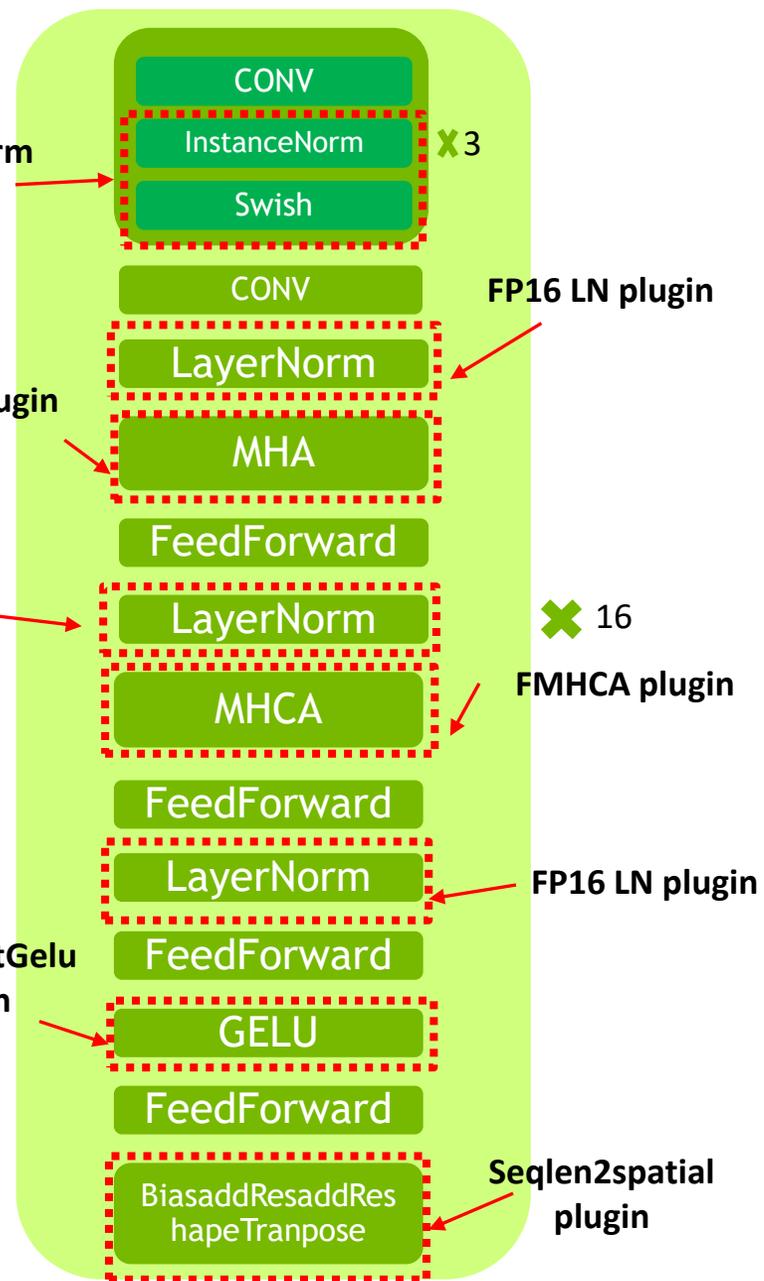
# U-NET

GroupNorm plugin

FMHA plugin

FP16 LN plugin

Flash SplitGelu plugin



# VAE

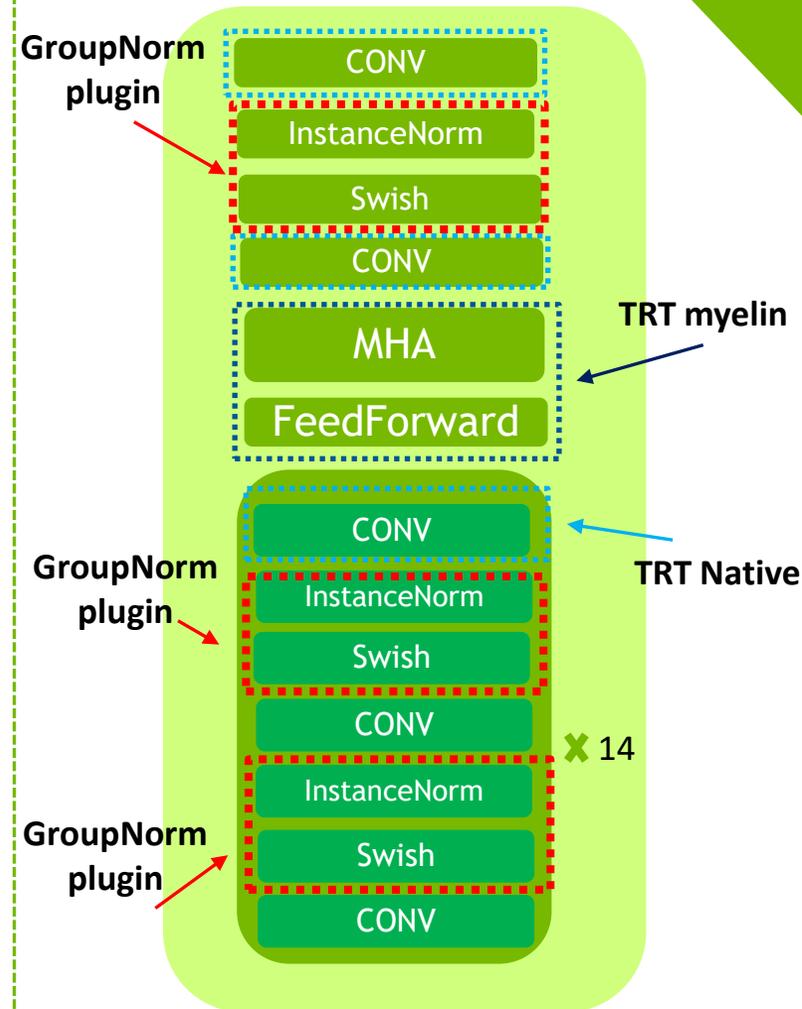
GroupNorm plugin

TRT myelin

GroupNorm plugin

GroupNorm plugin

TRT Native

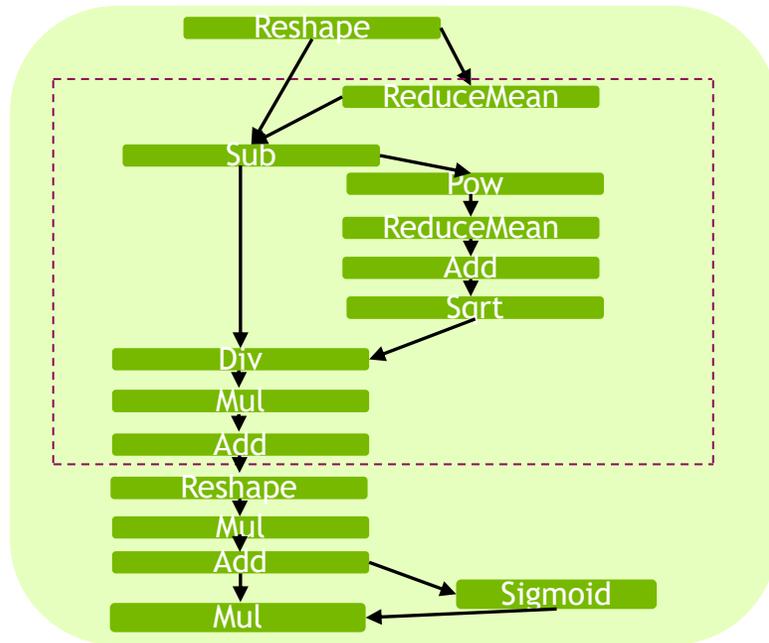




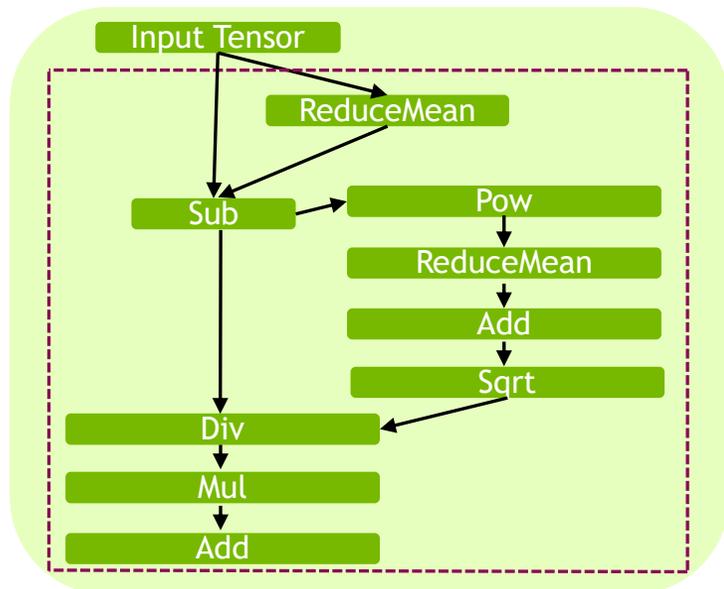
# Normalization

Breakup InstanceNormalization into primitive ops

InstanceNormalization



LayerNorm



Replace the whole pattern with GroupNorm Plugin

<https://github.com/NVIDIA/TensorRT/tree/main/plugin/groupNormPlugin>

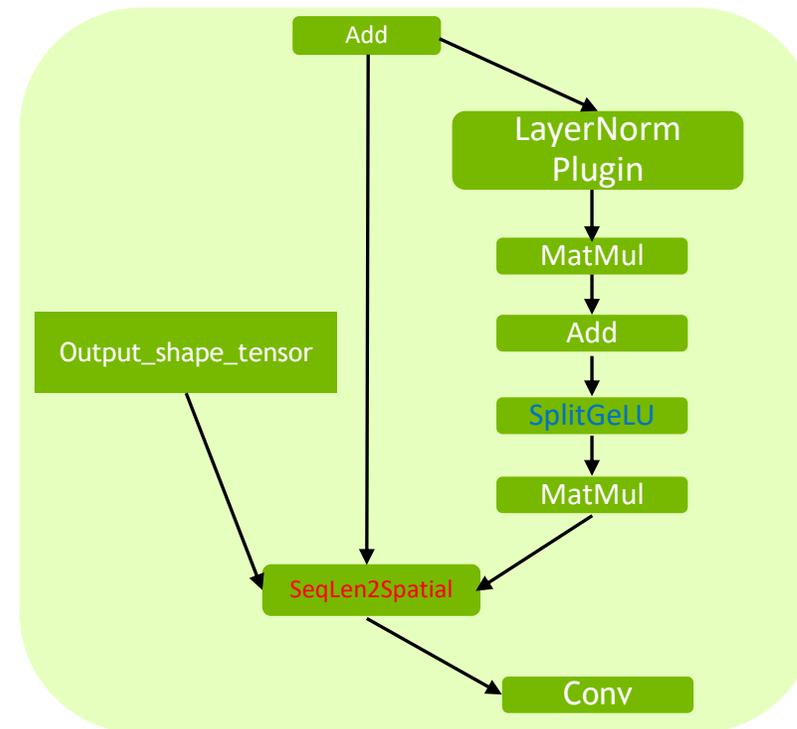
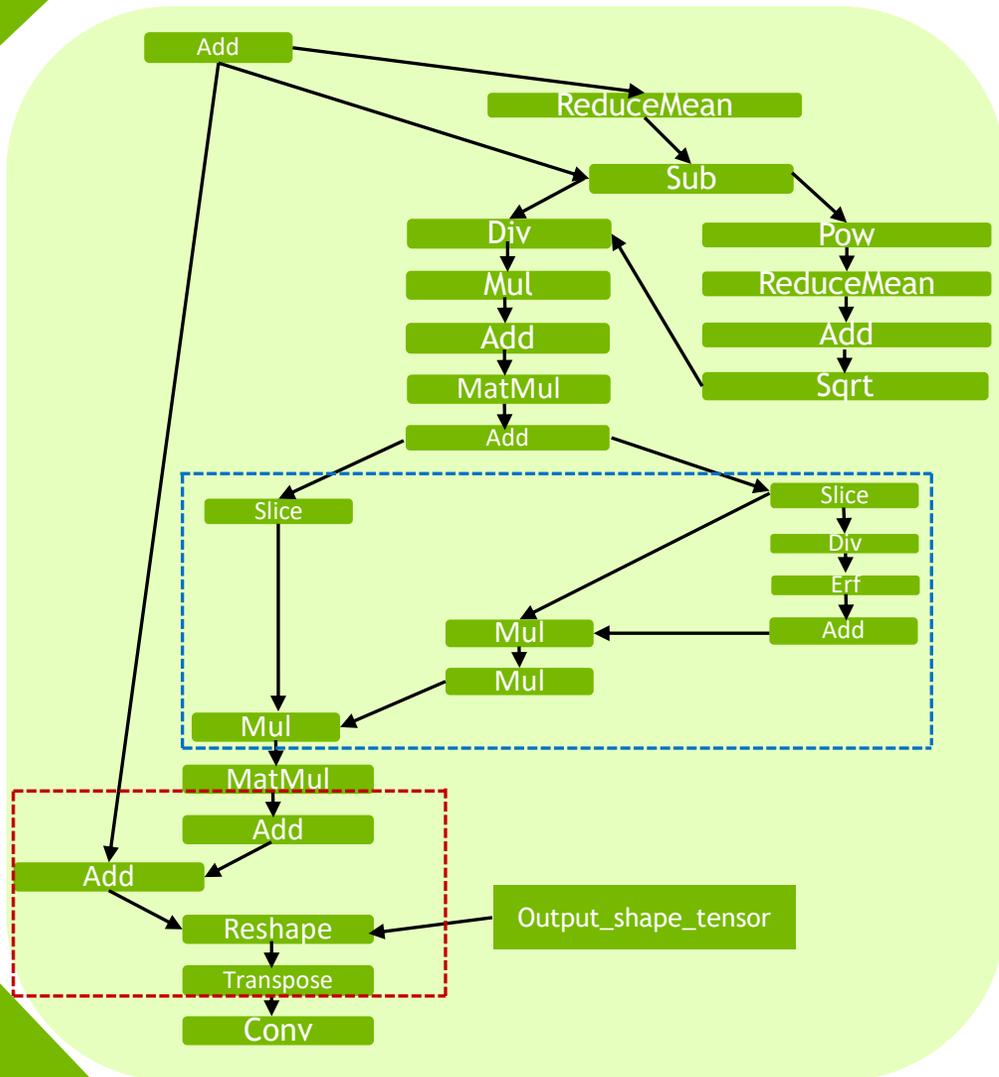
LayerNorm was also picked up by a plugin, since Myelin backend compiler was not used to optimize the MHA pattern

<https://github.com/NVIDIA/TensorRT/tree/main/plugin/layerNormPlugin>

```
#include "groupNormKernel.h"
#include <cuda/cub.cuh>
static inline __device__ __host__ float sigmoid(float x)
{
    return 1.F / (1.F + expf(-x));
}
struct GroupSums
{
    // Is it the 1st element of the group?
    int32_t flag;
    // The sum.
    float sum;
    // The sum of squares.
    float sumSq;
};
struct GroupSumsOp
{
    inline __device__ GroupSums operator()(GroupSums const& a, GroupSums const& b)
    {
        GroupSums dst;
        dst.sum = b.flag > b.sum ? (a.sum + b.sum);
        dst.sumSq = b.flag > b.sumSq ? (a.sumSq + b.sumSq);
        dst.flag = a.flag + b.flag;
        return dst;
    }
};
template <int32_t tTHREADS_PER_BLOCK>
__global__ void groupNorm@MHCKernel(GroupNorm@MHCParams params)
{
    // The object in charge of doing the sums for the different blocks.
    typedef cub::BlockScan<GroupSums, tTHREADS_PER_BLOCK> BlockScan;
    // Allocate shared memory for BlockScan.
    __shared__ typename BlockScan::TempStorage tempStorage;
    // Allocate shared memory for the groups. We could reduce the amount of shared
    // memory reserved.
    __shared__ float2 smem[tTHREADS_PER_BLOCK];
    // The instance in the batch.
    int32_t nI = blockIdx.z;
    // The channel loaded by that thread (2 channels per thread for F16x2).
    int32_t cI = blockIdx.x * params.cPerBlock + threadIdx.x * 2;
    // The first activation loaded by that block.
    int32_t hwBegin = blockIdx.y * params.hwPerBlock;
    // The last activation loaded by that block.
    int32_t hwEnd = min(hwBegin + params.hwPerBlock, params.hw);
```

```
#include "common/common.cuh"
#include "layerNormKernel.h"
template <typename T>
using kvp = cub::KeyValuePair<T>;
template <typename T>
struct mySum
{
    __host__ __device__ __forceinline__ kvp<T> operator()(kvp<T> const& a, kvp<T> const& b) const
    {
        return kvp<T>(a.key + b.key, a.value + b.value);
    }
};
template <typename T, typename OP_T, int32_t TPB>
__global__ void layerNorm@MHCKernel(
    int32_t const nHiddenDimension, T const* input, T const* gamma, T const* beta, T* output, float const epsilon)
{
    int32_t const index = blockIdx.x * nHiddenDimension + threadIdx.x;
    T const denominator = 1() / 1(nHiddenDimension);
    OP_T val = 0;
    kvp<OP_T> threadData(0, 0);
    if (threadIdx.x < nHiddenDimension)
    {
        val = input[index] * denominator;
        OP_T tmp0 = val * static_cast<OP_T>(denominator);
        OP_T tmp1 = val * tmp0;
        threadData = mySum<OP_T>()(threadData, kvp<OP_T>(tmp0, tmp1));
    }
    using MyReduce = cub::MyReduce<kvp<OP_T>, TPB>;
    __shared__ typename MyReduce::TempStorage temp;
    __shared__ OP_T mu, rSigma;
    auto const sumKV = MyReduce(temp).Reduce(threadData, mySum<OP_T>());
    if (threadIdx.x == 0)
    {
        mu = sumKV.key;
        rSigma = rsqrt(sumKV.value - mu * mu + static_cast<OP_T>(epsilon));
    }
```

# Other Small Fusions(SplitGeLU+SeqLen2Spatial)



\*Additionally, we also employ **SplitGelu plugin** to fuse ops around Erf and **seq2spatial plugin** to fuse the biasadd+residualadd+reshape+transpose pattern



# DemoDiffusion Benchmark

Optimization Level	UNETx50
OOTB(out of the box)	2288.14ms
Myelin (RemoveCast + RemoveInstanceNorm)	1282.32ms
Above + fMHA +fMHCA	1047.26ms
Above + LayerNorm	979.72ms
Above + GroupNorm	874.51ms
Above + SplitGeLU + Seq2Spatial	856.47ms
Above + Graph Optimization	805.97ms
Above + Preview Feature 0805	792.26ms

Optimization Level	CLIP
OOTB(out of the box)	3.01ms
Myelin (RemoveCast + RemoveInstanceNorm)	2.96ms
Above + LayerNorm	4.38ms

Optimization Level	VAE
OOTB(out of the box)	20.73ms
Myelin (RemoveCast + RemoveInstanceNorm)	20.97ms
Above + GroupNorm	18.54ms

TensorRT8.5.2.2, batch=1, image  
size=512x512,  
GPU=A100-80Gib-PCIE

## Pipeline Total Time

815.18ms

Support dynamic shapes image generation up to  
1024x1024

TensorRT8.5.2, batch=1, image size=512x512,GPU=L40

Optimization Level	CLIP	UNETx50	VAE	Pipeline
ALL Plugins	2.42ms	973.56ms	23.98ms	999.98ms

TensorRT8.5.1, batch=1, image size=512x512,GPU=T4

Optimization Level	CLIP	UNETx50	VAE	Pipeline
ALL Plugins	9.70ms	4536.414ms	104.944ms	4651.088ms

TensorRT8.5.1, batch=1, image size=512x512,GPU=A10

Optimization Level	CLIP	UNETx50	VAE	Pipeline
ALL Plugins	5.14ms	1892.144ms	40.180ms	1937.497ms

TensorRT8.5.1, batch=1, image size=512x512,GPU=A30

Optimization Level	CLIP	UNETx50	VAE	Pipeline
ALL Plugins	4.97ms	1368.08ms	30.34ms	1403.43ms