

Real data labeling → Text to 4D hybrid dataset ← Text to 4D AI generation

Text: An adventurer is riding dinosaur through the rainforest.

Text Extended Template

A Subject Details:

An adventurer wearing a hat and backpack sitting on a brown dinosaur in a vibrant rainforest with sunlight filtering through dense foliage and a stream below.

B Action Pattern:

Riding the dinosaur across a flowing stream through the lush jungle terrain with dynamic movement and splashing water.

T5 Text Encoder

1- Global text alignment

β_1

Object Text Encoder

α_1

β_2

Motion Text Encoder

α_2

Quarter Division
(Data Matrix $N \times N$)

VAE Encoder

Timestep

Noise

Z_n

Self Attention Layer

View cross attention

Frame cross attention

Z_{n-1}

3- Guided Diffusion1

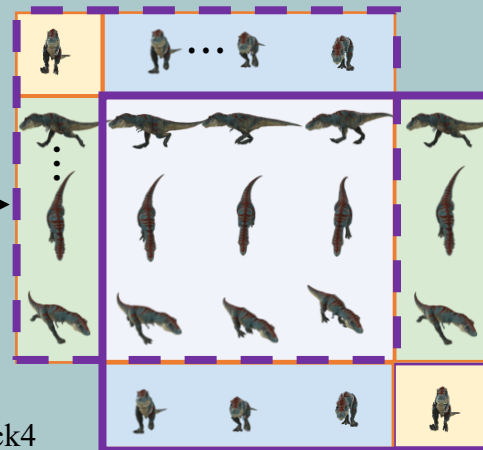
(1) I-Swin nine-grid division

Conditional Guidance1



Nine-grid Matrix Division

(2) Swin Shifting



Block4

(3) Slicing window sampling



(4) Conditional Guidance2



View cross attention

$1/2N \times 1/2N$



Frame cross attention

$1/2N \times 1/2N$

Sliding Slice Sampling

AttentionMask1

AttentionMask2

4- Guided Diffusion2

2- Fine grained text alignment