

RecipeSD: Injecting Recipe into Food Image Synthesis with Stable Diffusion



Jing Yang

Junwen Chen

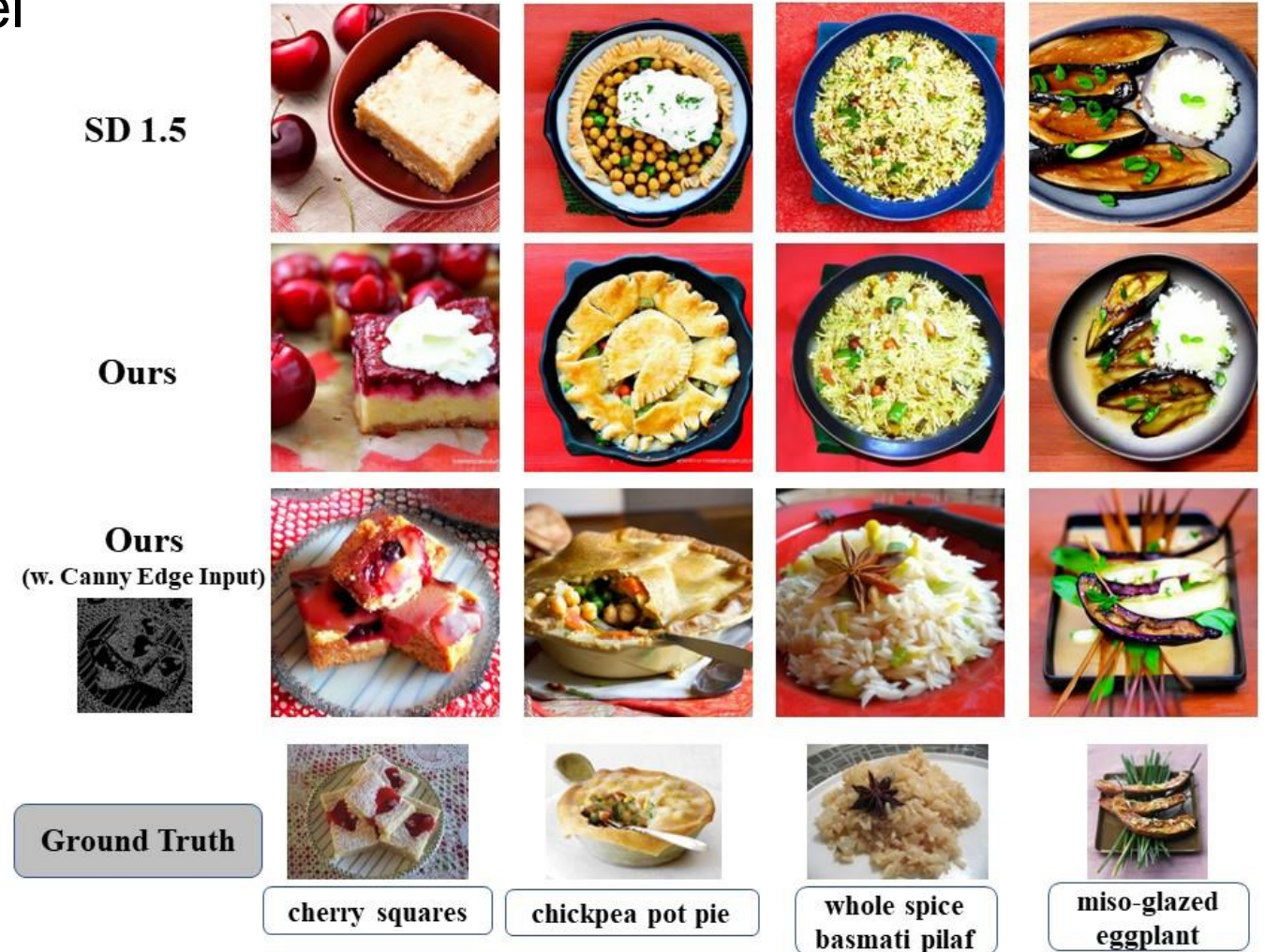
Keiji Yanai



The University of Electro-Communications

□ RecipeSD

- Transfer cross-modal retrieval model knowledge to SD
- Introduce recipe text embeddings
- Incorporate with ControlNets



□ Modality

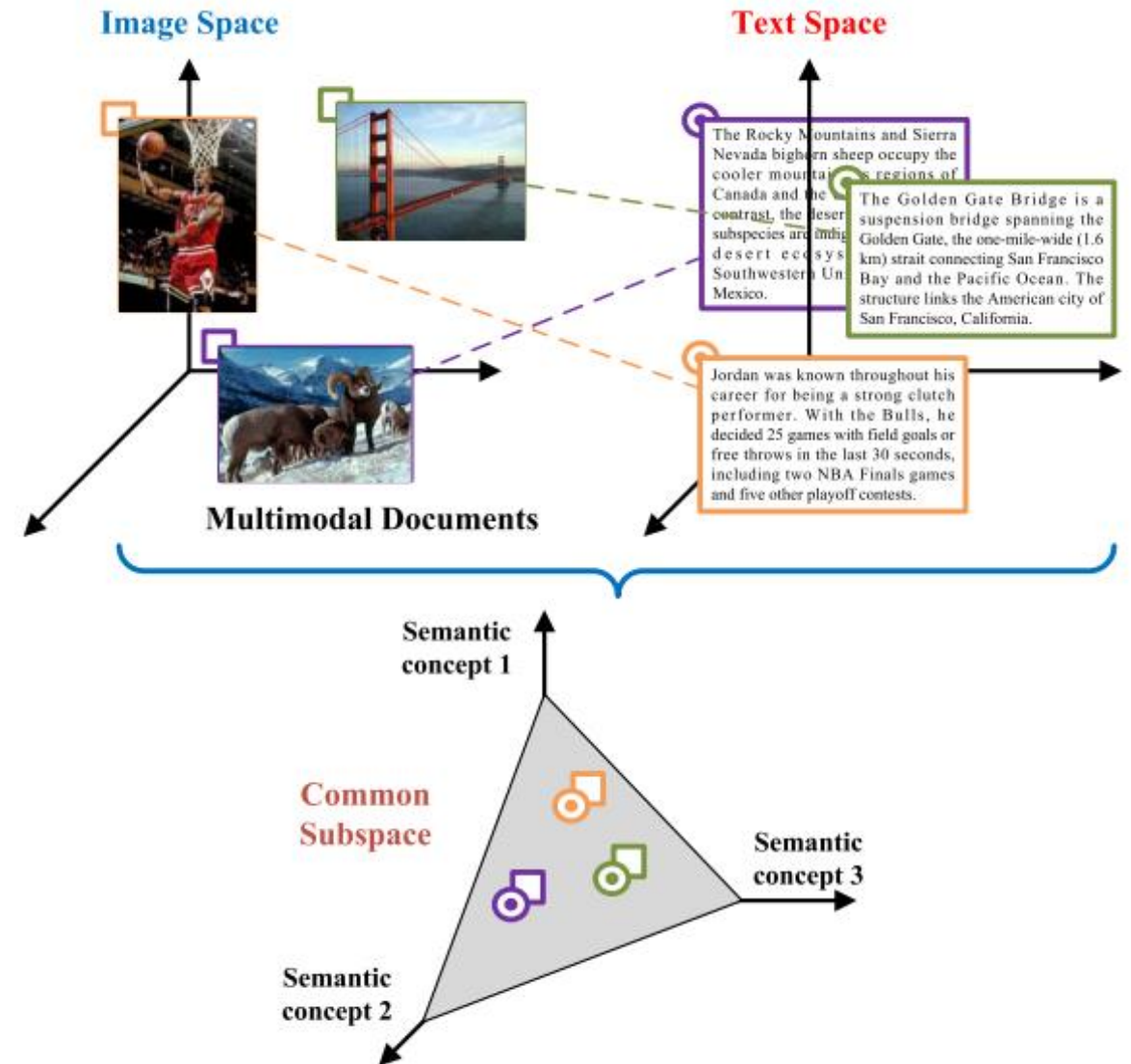
- Text, image, audio, video...

□ Cross-modal image-text retrieval

- Build the connection is difficult
➡ The gap between modalities

□ Solution

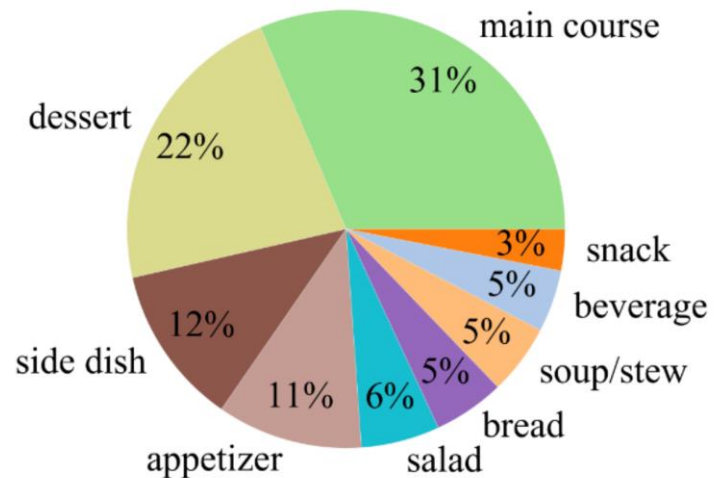
- Embeddings & Distance Learning
- A large number of data pairs



Recipe Retrieval and Dataset

□ Recipe1M

- One of the applications of cross-modal retrieval
- 1 million pairs of recipe images and recipe texts



Query Image



Retrieved Recipe

Ingredients	Instructions
butter	1. Melt 1 tablespoon butter with 1/2 tablespoon olive oil in saucepan over medium heat.
olive oil	2. Add onions and saute, stirring every few minutes, until they are caramelized, about 15-20 minutes.
sweet onions	...
portabella mushrooms	3. (If soup is too thick, thin with a little more hot broth).
celery	4. Season to suit your taste with salt and freshly-cracked black pepper.
carrot	5. Serve in deep bowls, garnished with a sprinkle of minced, fresh parsley.
garlic cloves	
...	

Query Recipe

Ingredients	Instructions
hamburger	1. Cook hamburger until done and drain off the fat.
rigatoni pasta	2. Add mushrooms and onion and fry until translucent.
Ragu pizza sauce	3. Add pepperoni.
mushrooms	4. ...
onion	5. Lay noodles on top of hamburger mix in crockpot.
pepperoni	6. Turn crock on low and leave 4-5 hours.
mozzarella cheese	7. Pour over the remainder of pizza sauce over the noodles.
	8. Top with the cheese.



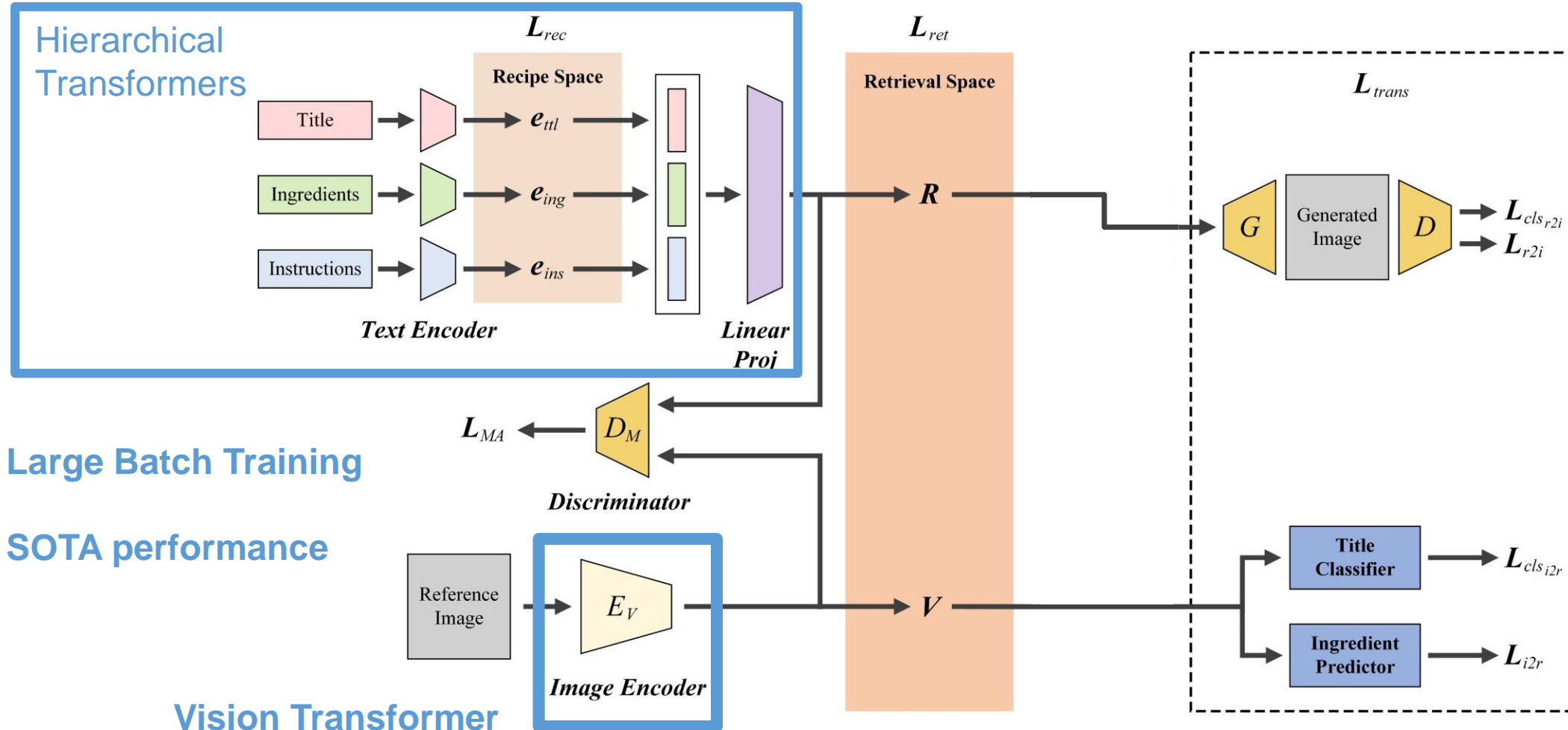
Retrieved Image

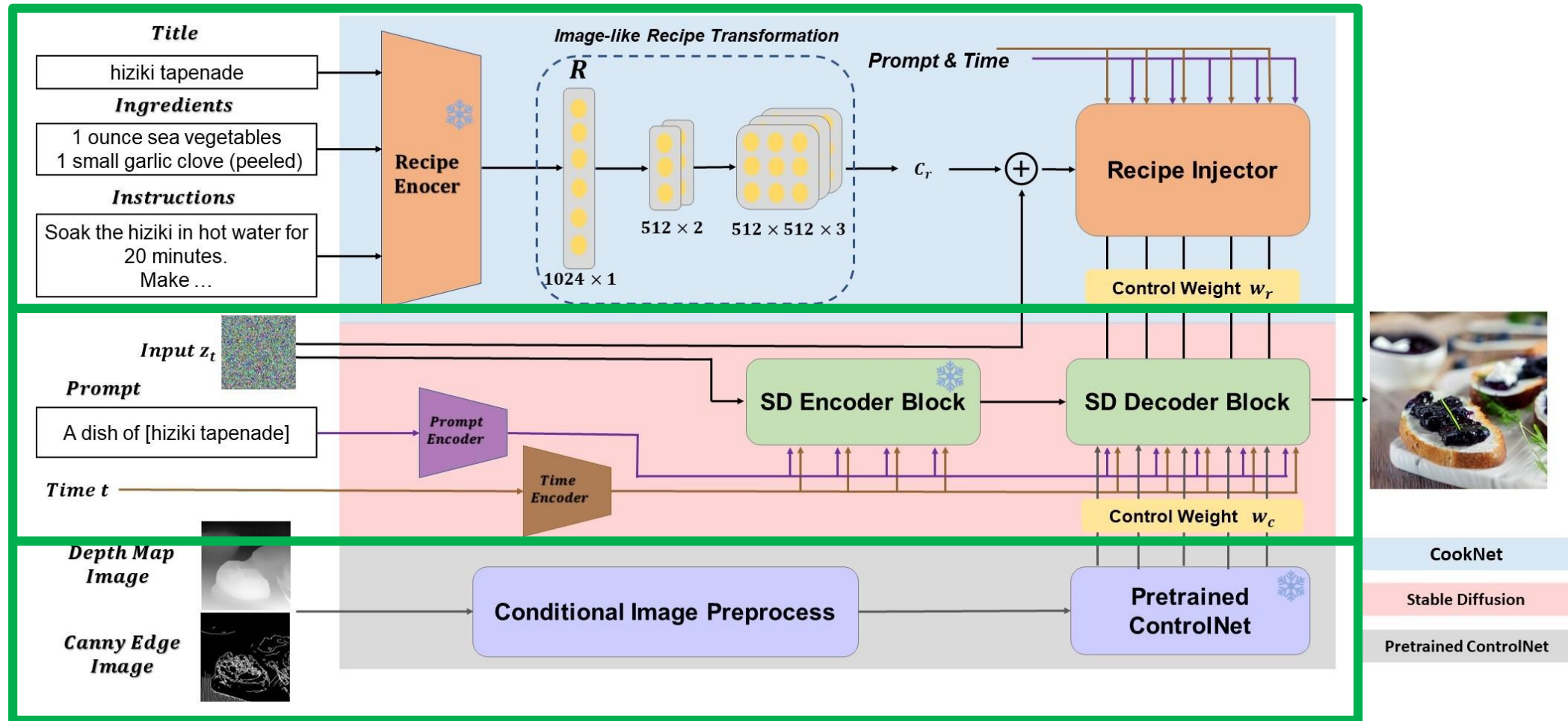


Recipe Retrieval Method

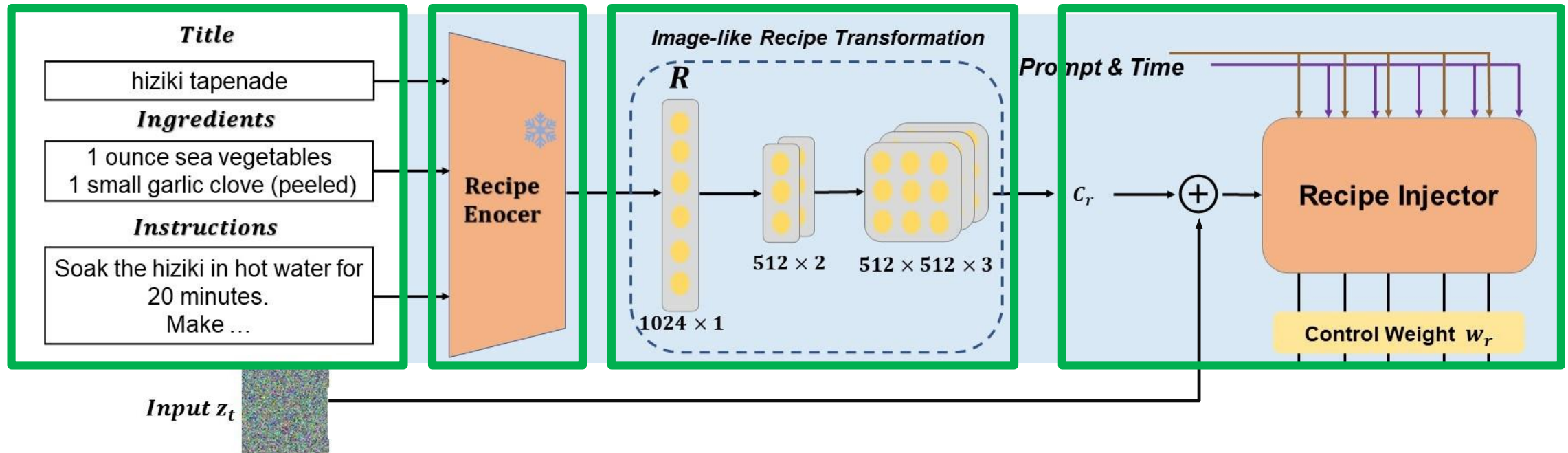
□ TNLBT

- Improving the representation capability of the recipe embeddings

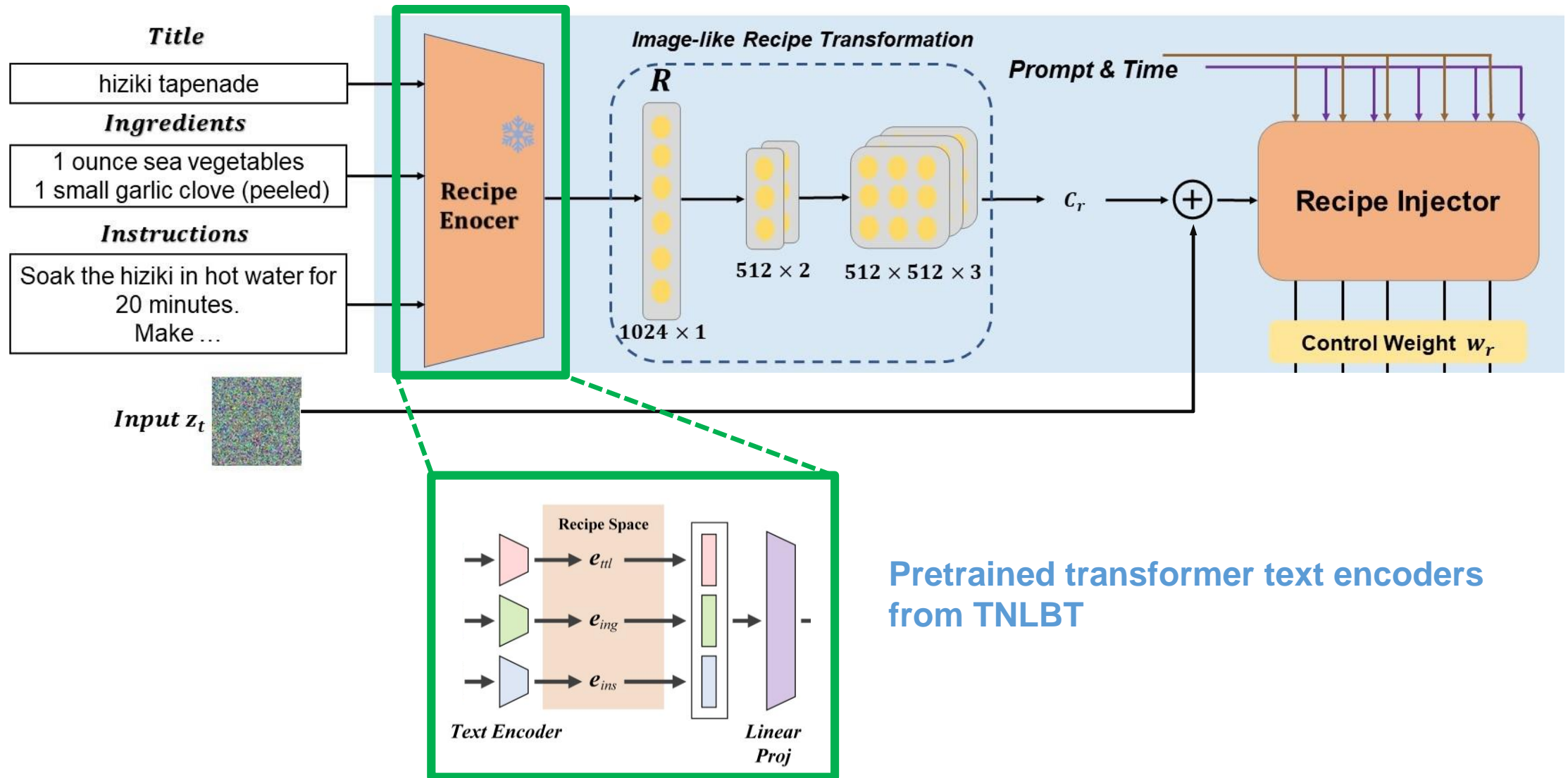




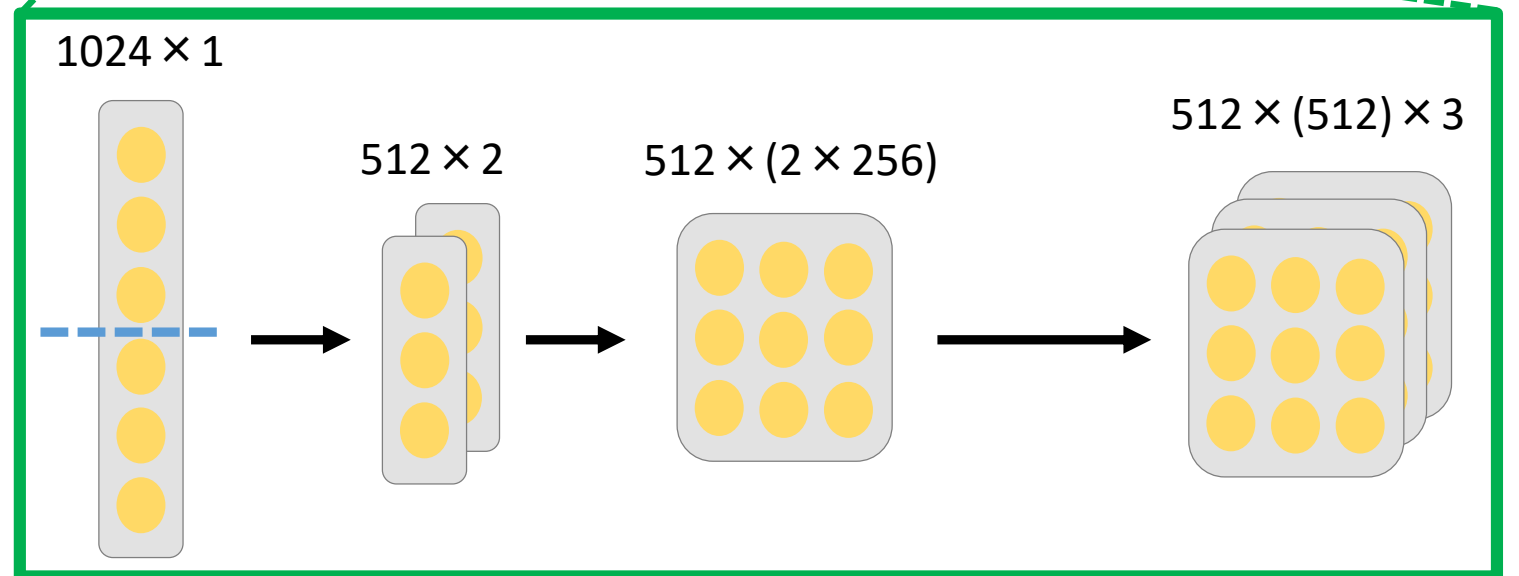
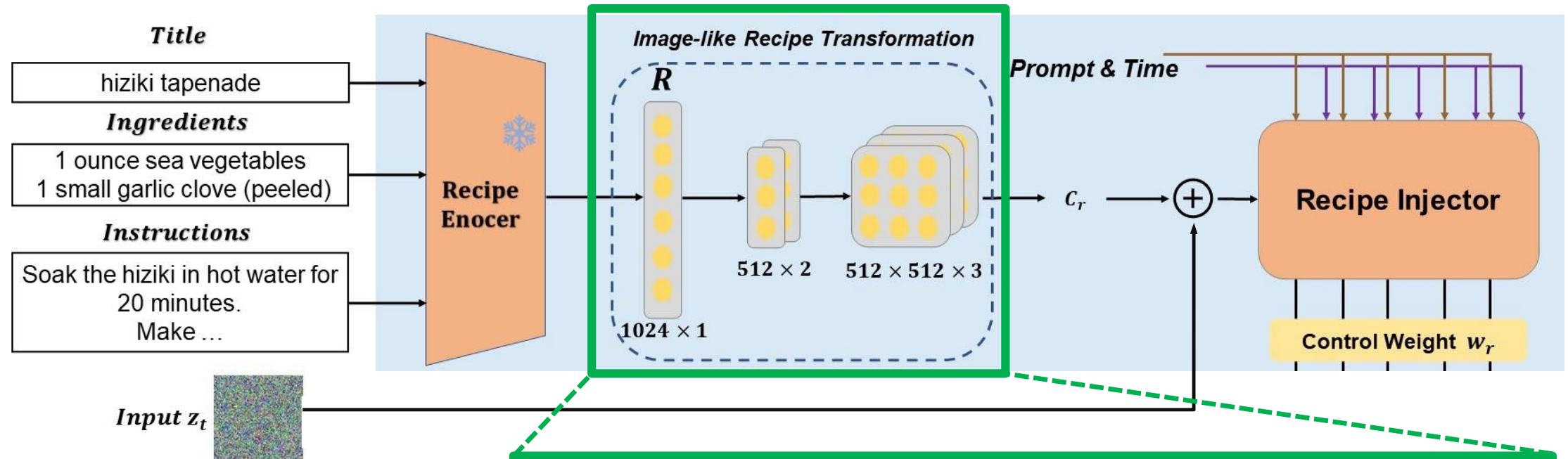
- Our RecipeSD consists of three parts: Stable Diffusion, CookNet, and ControlNet
- Image-like Recipe Transformation (IRT) transforms the recipe text embeddings to image-like embeddings



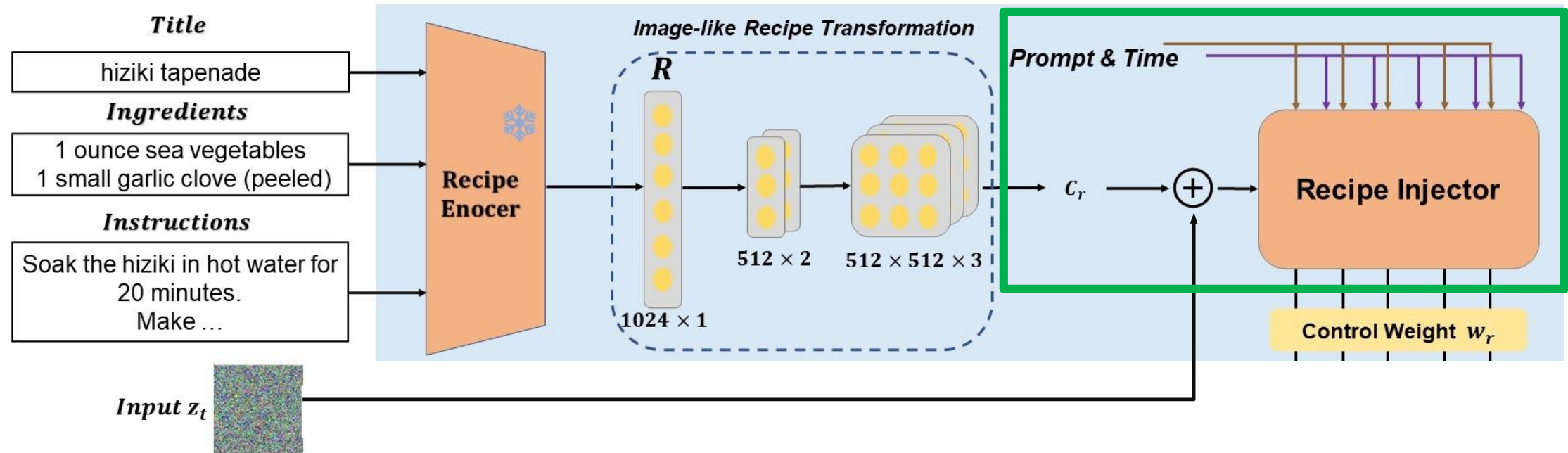
- Distinct control between different detailed recipe texts
- Image generation based on recipe image embeddings
- Reconstruct image information from recipe text
- Accurately reproducing the specified ingredients



- Prior knowledge from cross-modal retrieval task



- Edit-free recipe embedding transformation
- Directly save text information to pixel



- Our Recipe Injector incorporate reconstruct text information from the beginning
- Accurately reproducing the specified ingredients



Hiziki tapenade



Berries romanoff



Dianne's lemon-feta quinoa salad



Stir-fry pork with ginger



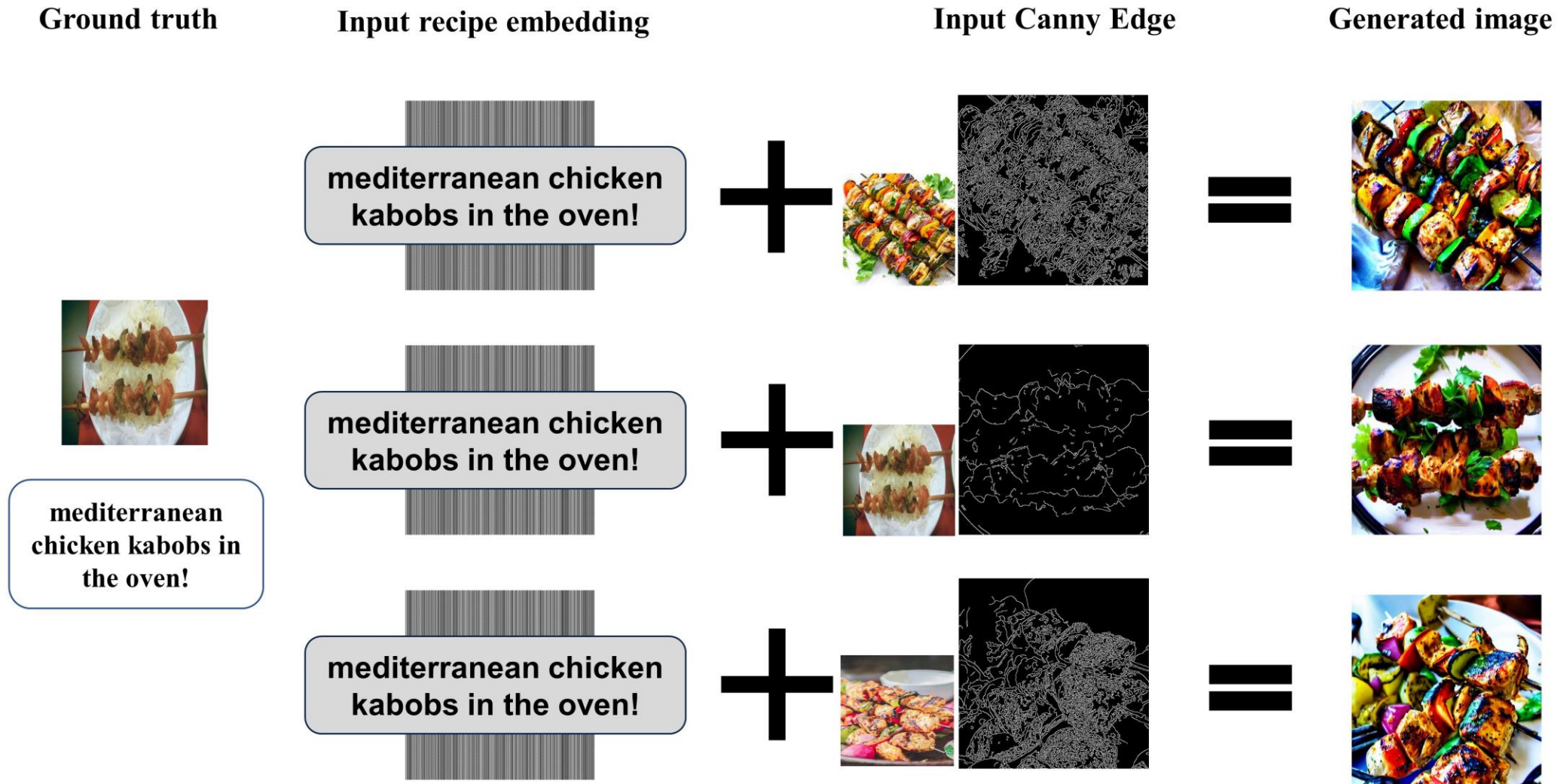
Ground Truth

0

recipe control weight increasing

2.0

With the detailed recipe information from recipe text encoders and IRT in our CookNet, we increasing the control weight to improve the quality of generated images



Recipe Embedding + Canny Edge

Ground truth

Input recipe embedding

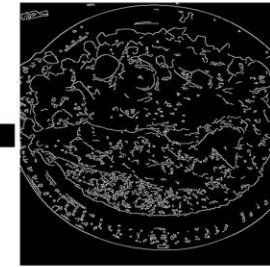
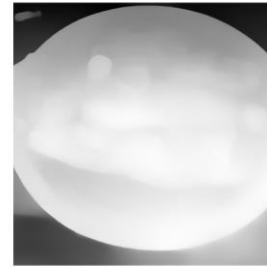
Input Depth image

Input Canny Edge

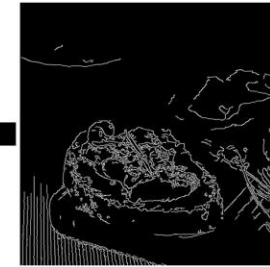
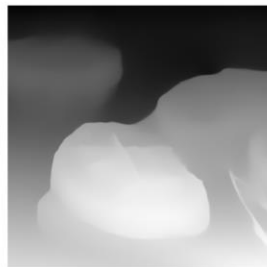
Generated image



**citrus steamed trout
with quinoa pilaf**



**blueberry ginger jam
& goat cheese crostini**



**carrot cookies with
cinnamon cream
cheese icing**



Recipe Embedding + Depth Image + Canny Edge

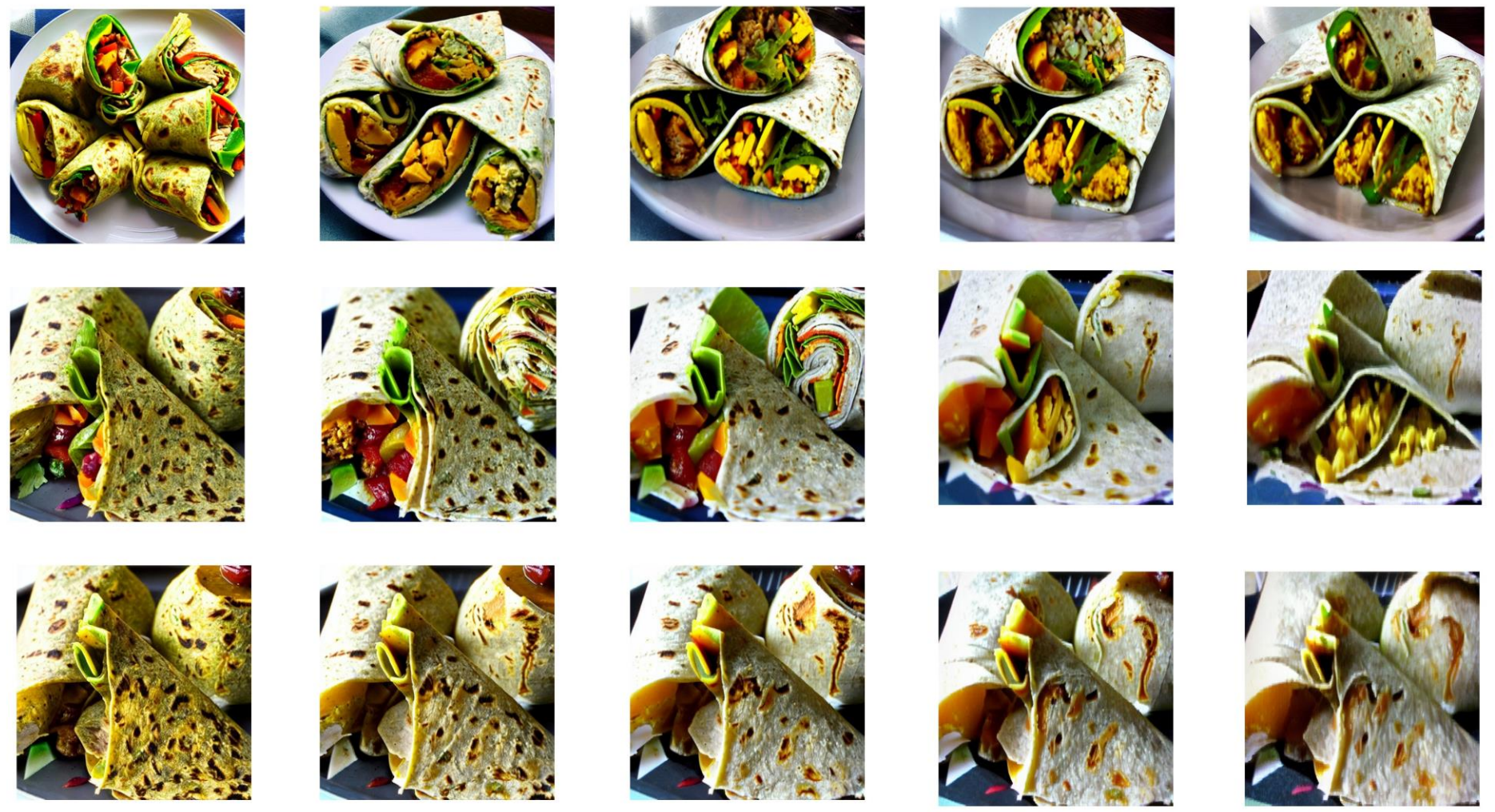
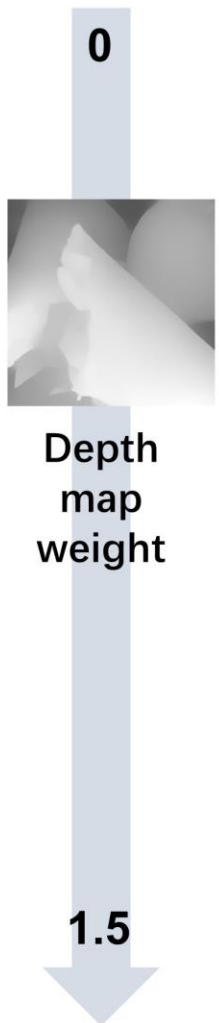


curried turkey wraps

0

recipe embedding weight

2.0



- We introduce RecipeSD for synthesizing food images by injecting recipe information into the Stable Diffusion model
- We proposed CookNet with task-specific recipe encoders to align the generated images with the corresponding recipe texts
- We demonstrated that our approach can be further enhanced by incorporating other ControlNets

