

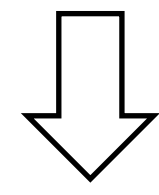
## Motivation for Better Knowledge Distillation

### Representation Gap in KD

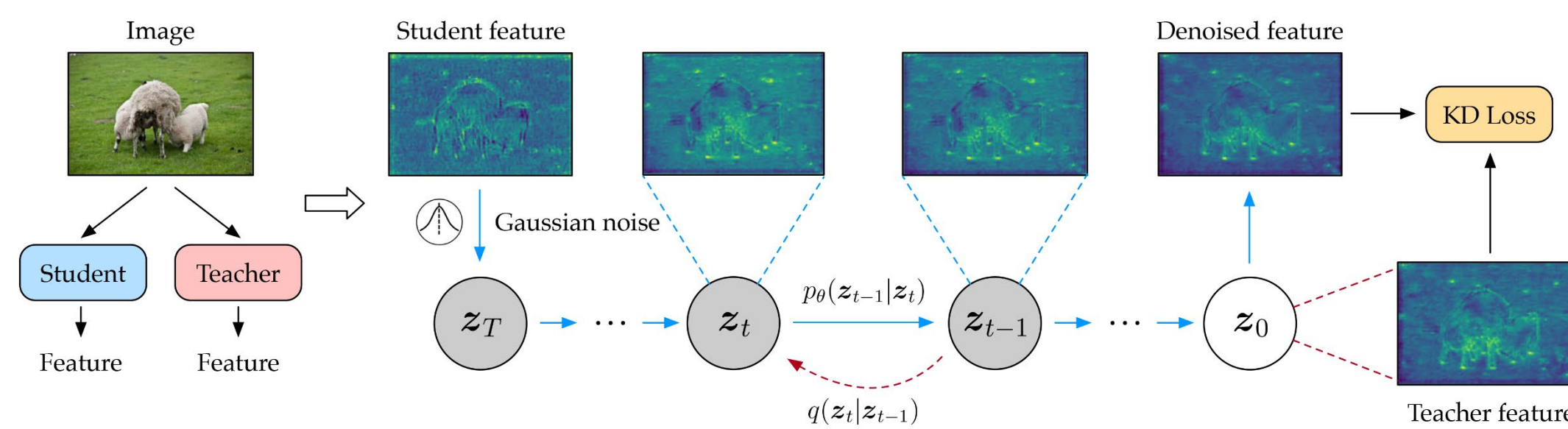
- Challenge of the teacher-student model capacity gap.
- Existing methods are often complex and task-specific.

### Noise in Distillation Features

- Student features are noisier due to the limited capacity.
- The noise leads to suboptimal distillation and performance.



## Our DiffKD Approach



- A novel method using diffusion models for denoising the student features.
- Distillation on denoised student features with simple losses such as MSE.

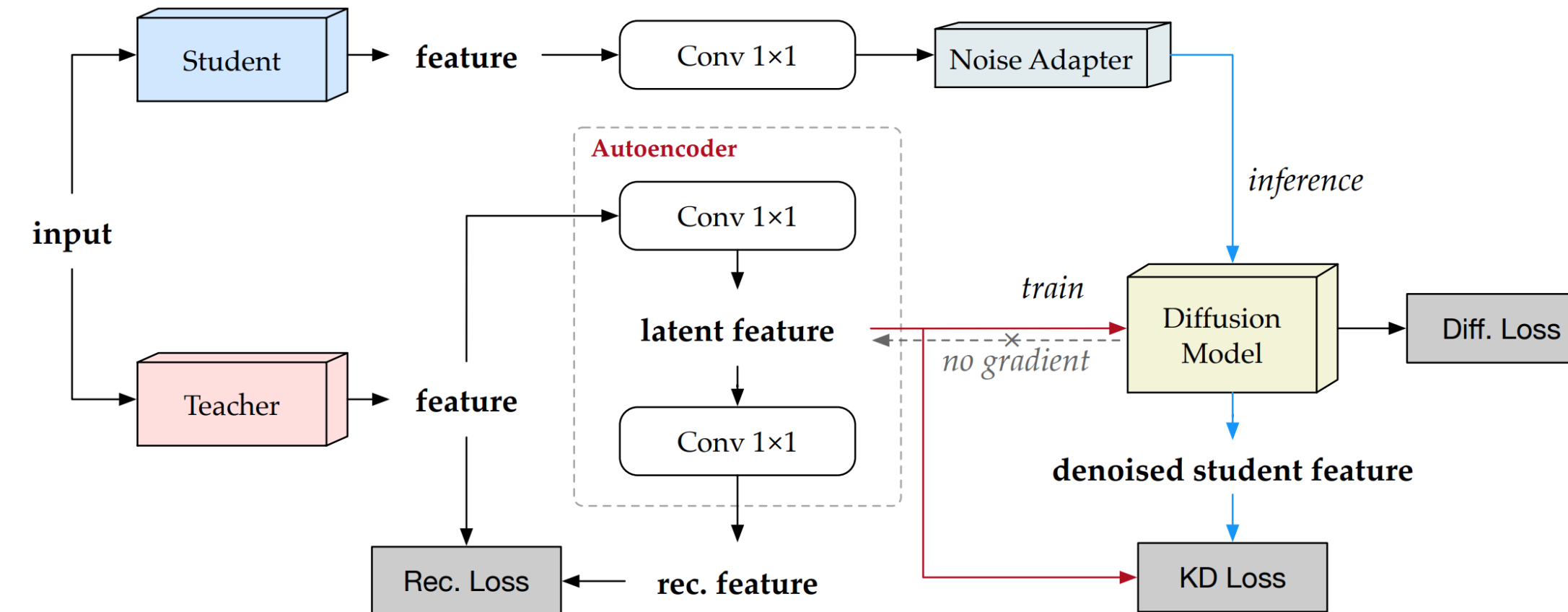
### Innovations

- Lightweight diffusion model with linear autoencoder.
- Adaptive noise matching for precise denoising.

### Effectiveness

- Applicable to various feature types.
- Superior performance in multiple tasks and settings.

## Method

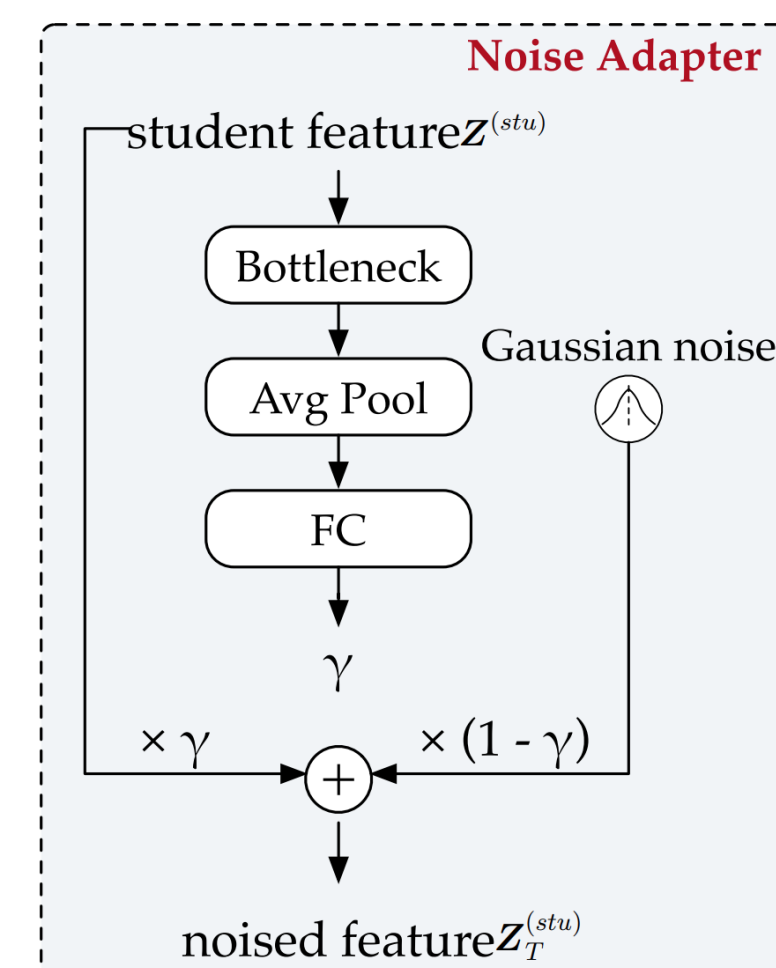
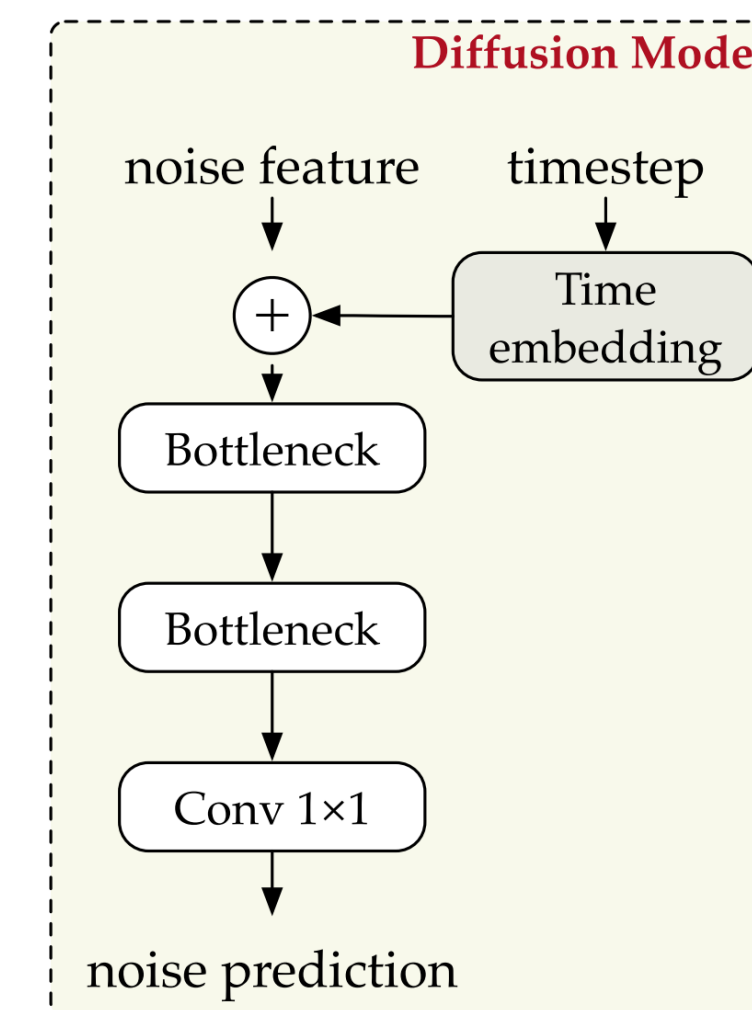


### Simultaneous Optimization with

- *Task loss* for training the student.
- *KD loss* for training the student & noise adapter.
- *Diffusion loss* for training the diffusion model.
- *Reconstruction loss* for training the linear autoencoder.

### Diffusion Model

- A lightweight model with ResNet Bottleneck blocks.
- Trained with teacher features.
- Leveraged for denoising student features.



### Noise Adapter

- Addresses the challenge of inexact noisy levels in student features.
- Measures the noisy level of feature.
- Complements additional Gaussian noise to feature to match the noisy level.

## Experiments

### ImageNet

Student (teacher)	Tea.	Stu.	KD [13]	Review [6]	DKD [50]	DIST [16]	MSE	DiffKD	DiffKD <sup>†</sup>	
R18 (R34)	Top-1	73.31	69.76	70.66	71.61	71.70	72.07	70.58	72.22	<b>72.49</b>
	Top-5	91.42	89.08	89.88	90.51	90.41	90.42	89.95	90.64	<b>90.71</b>
MBV1 (R50)	Top-1	76.16	70.13	70.68	72.56	72.05	73.24	72.39	73.62	<b>73.78</b>
	Top-5	92.86	89.49	90.30	91.00	91.05	91.12	90.74	91.34	<b>91.48</b>

### ImageNet with Stronger Teachers

Teacher	Student	Top-1 ACC (%)						
		Tea.	Stu.	KD [13]	RKD [30]	SRRL [46]	DIST [16]	DiffKD
ResNet-50	ResNet-34		76.8	77.2	76.6	76.7	77.8	<b>78.1</b>
	MobileNetV2	80.1	73.6	71.7	73.1	69.2	74.4	<b>74.9</b>
	EfficientNet-B0		78.0	77.4	77.5	77.3	78.6	<b>78.8</b>
Swin-L	ResNet-50	86.3	78.5	80.0	78.9	78.6	80.2	<b>80.5</b>
	Swin-T		81.3	81.5	81.2	81.5	82.3	<b>82.5</b>

### COCO

Method	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<i>Two-stage detectors</i>						
T: Faster RCNN-R101	39.8	60.1	43.3	22.5	43.6	52.8
S: Faster RCNN-R50	38.4	59.0	42.0	21.5	42.1	50.3
DiffKD	40.6	60.9	43.9	<b>23.0</b>	44.5	<b>54.0</b>
T: CM RCNN-X101	45.6	64.1	49.7	26.2	49.6	60.0
S: Faster RCNN-R50	38.4	59.0	42.0	21.5	42.1	50.3
DiffKD	42.2	62.8	46.0	<b>24.2</b>	46.6	55.3

### Cityscapes

Method	Params (M)	FLOPs (G)	mIoU (%)	
			Val	Test
T: DeepLabV3-R101	61.1	2371.7	78.07	77.46
S: DeepLabV3-R18	13.6	572.0	74.21	73.45
DiffKD	13.6	572.0	<b>77.78</b>	<b>76.24</b>
T: DeepLabV3-R101	61.1	2371.7	78.07	77.46
S: PSPNet-R18	12.9	507.4	72.55	72.29
DiffKD	12.9	507.4	<b>75.83</b>	<b>75.61</b>

## Visualizations

