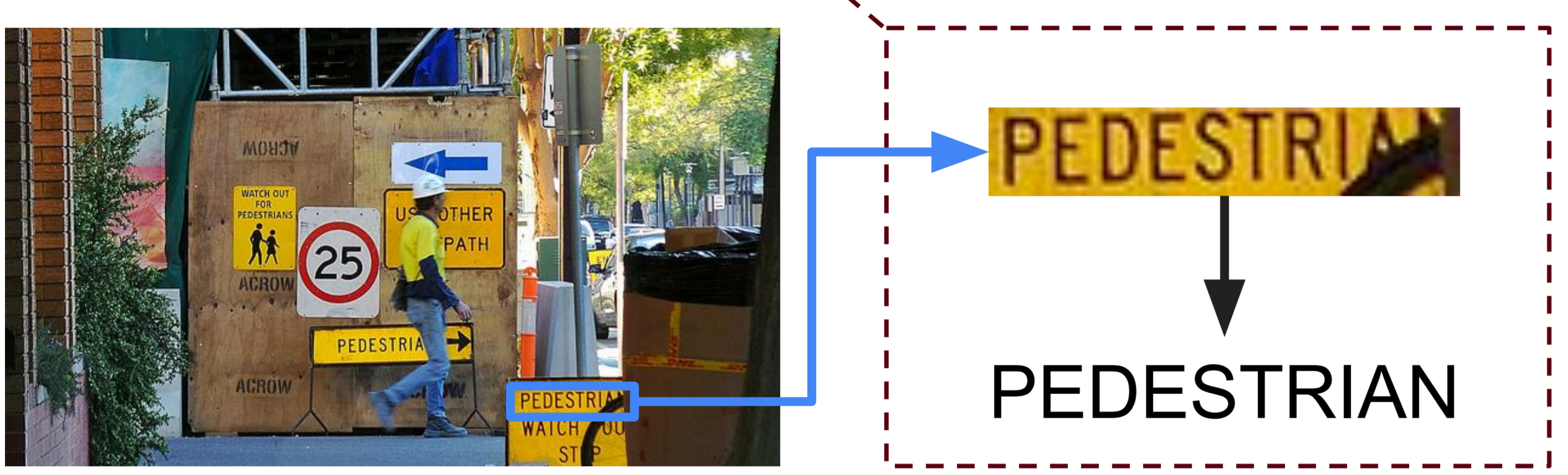


## Problem Text recognition in natural scenes



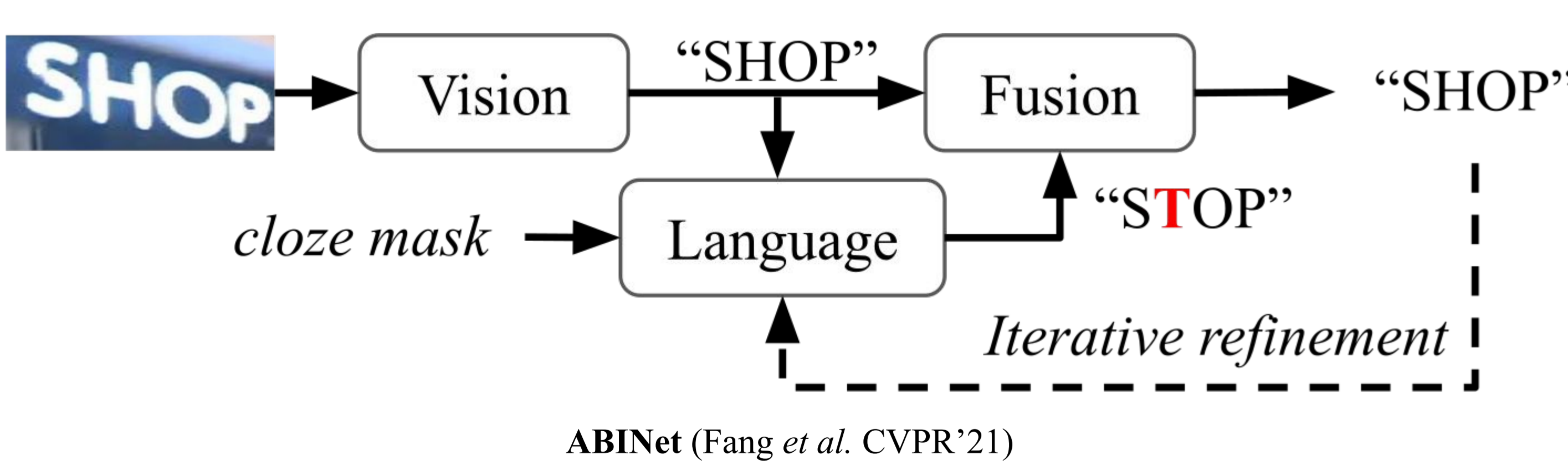
## Language context can resolve visual ambiguity



## But context-aware STR methods are typically limited to:

- Monotonic decoding (one character at a time)
- Unidirectional context (usually left-to-right)

## Recent two-stage approaches are inefficient



## Language Model accounts for 35% of the parameters but is:

- Underutilized (uses only 13.65% of total FLOPS)
- Error-prone due to lack of visual context (50.44% word accuracy with ground truth label as input)

## Image is the primary input signal in STR, not language context

→ The Language Model should also consider visual context

## Key Idea Unify models using AR ensemble

Ensemble of AR models (PARSeq model)

$$P(y|\mathbf{x})_{[1,2,3]} = P(y_1|\mathbf{x})P(y_2|y_1, \mathbf{x})P(y_3|y_1, y_2, \mathbf{x})$$

$$P(y|\mathbf{x})_{[3,2,1]} = P(y_3|\mathbf{x})P(y_2|y_3, \mathbf{x})P(y_1|y_2, y_3, \mathbf{x})$$

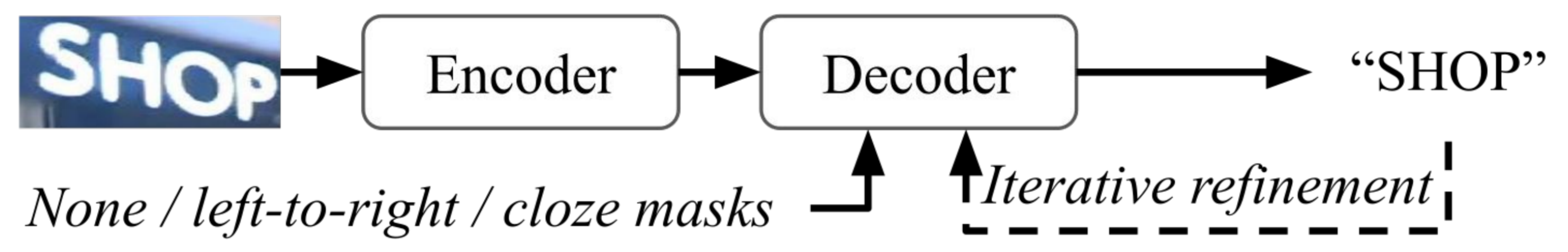
$$P(y|\mathbf{x})_{[1,3,2]} = P(y_1|\mathbf{x})P(y_3|y_1, \mathbf{x})P(y_2|y_1, y_3, \mathbf{x})$$

$$P(y|\mathbf{x})_{[2,3,1]} = P(y_2|\mathbf{x})P(y_3|y_2, \mathbf{x})P(y_1|y_2, y_3, \mathbf{x})$$

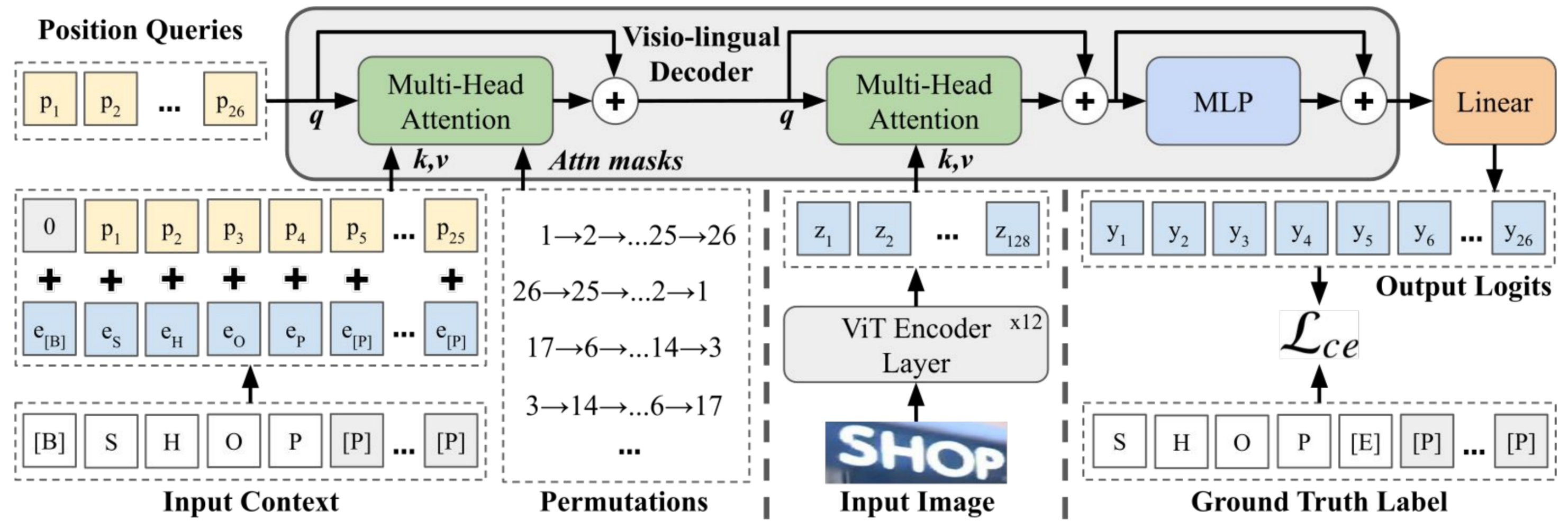
$$\prod_{t=1}^T P(y_t|y_{<t}, \mathbf{x}) \quad \prod_{t=1}^T P(y_t|\mathbf{x}) \quad \prod_{t=1}^T P(y_t|y_{\neq t}, \mathbf{x})$$

Context-aware AR model      Context-free NAR model      Bidirectional Iterative Refinement model

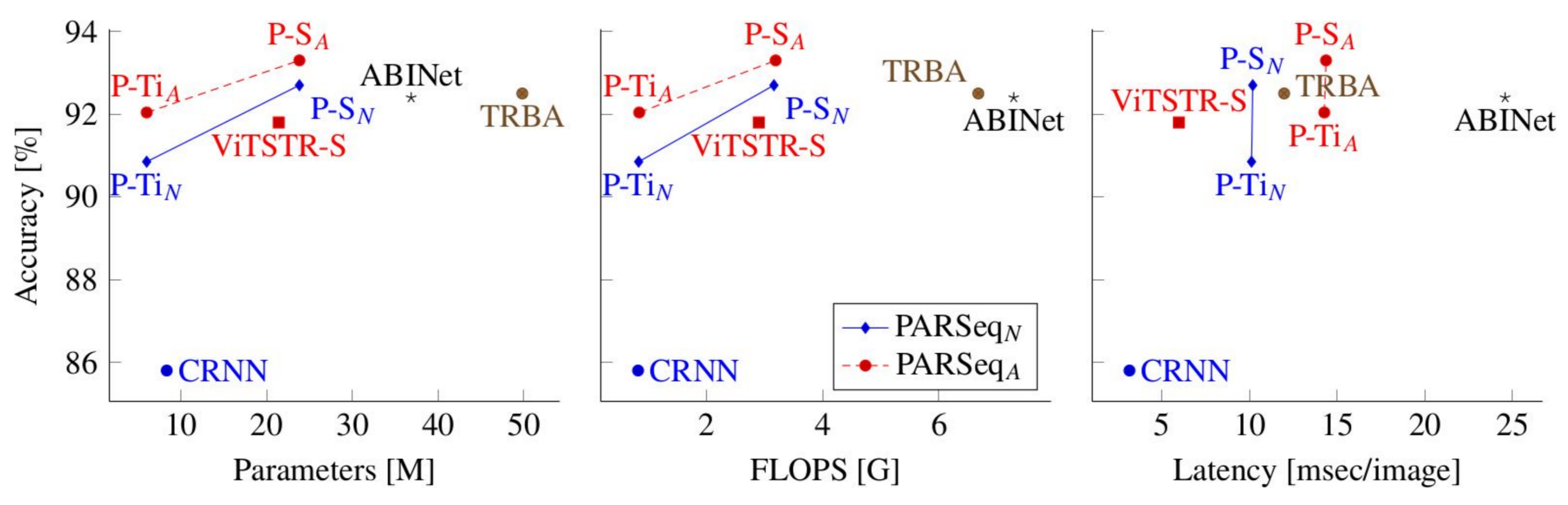
## PARSeq: A unified model for STR



## Learning PARSeq with Permutation Language Modeling:



## PARSeq is flexible, accurate, and efficient

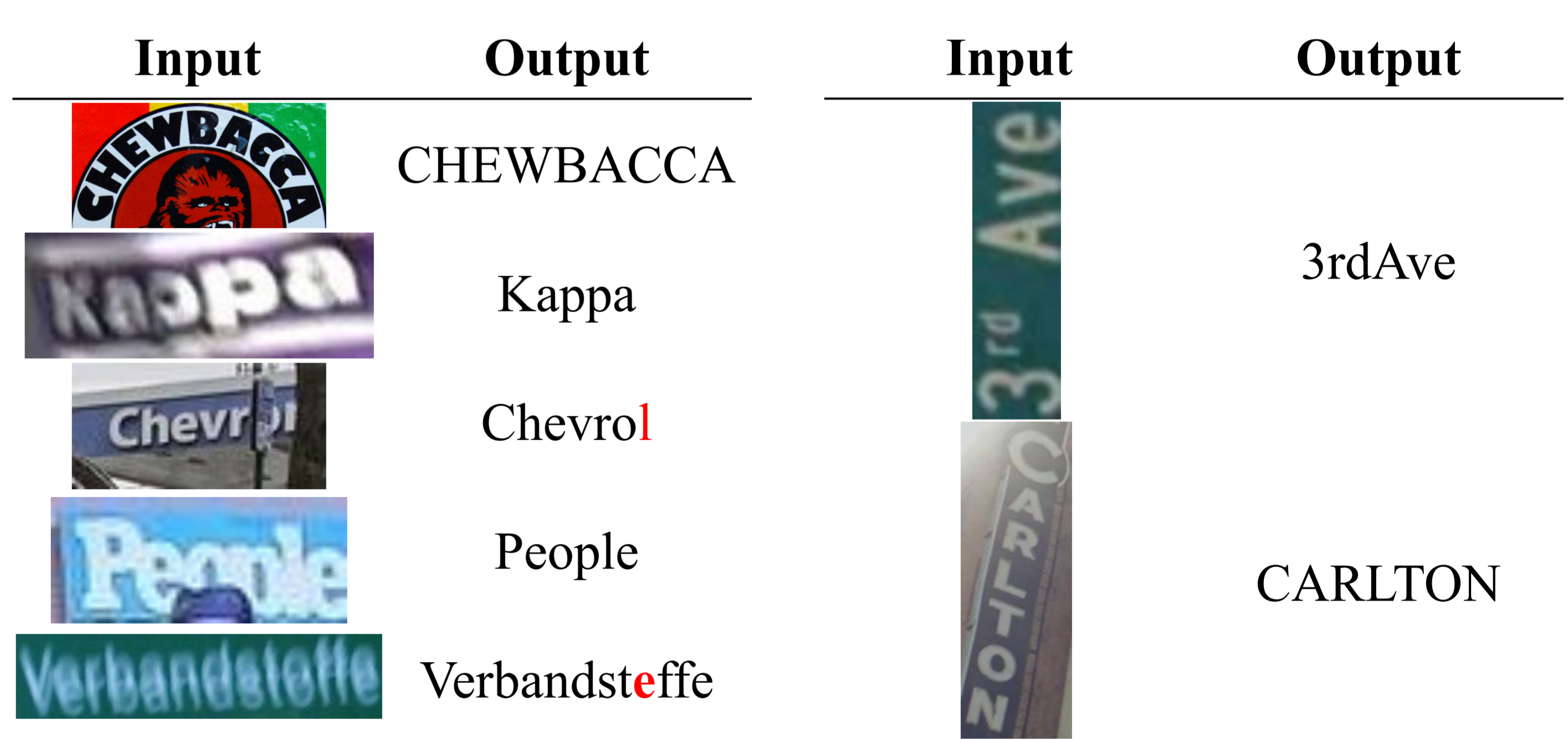


- PARSeq<sub>N</sub> for constant-time inference (non-autoregressive)
- PARSeq<sub>A</sub> for more accurate results (autoregressive)

## Results SOTA in STR benchmarks

Previous Work	Method	Conf.	Train data	36-char word acc.	
				7,248 samples	7,672 samples
Previous Work	SRN	CVPR'20	MJ,ST	90.4	–
	TextScanner	AAAI'20	MJ,ST+	–	91.0
	Bhunja et al.	ICCV'21	MJ,ST	–	90.9
	VisionLAN	ICCV'21	MJ,ST	91.2	–
	PREN2D	CVPR'21	MJ,ST	91.5	–
	ABINet	CVPR'21	MJ,ST+	92.7	–
Ours	PARSeq <sub>N</sub>		MJ,ST	92.0±0.2	90.7±0.2
	PARSeq <sub>A</sub>		MJ,ST	<b>93.2±0.2</b>	<b>91.9±0.2</b>
	PARSeq <sub>N</sub>		real	95.7±0.1	95.2±0.1
	PARSeq <sub>A</sub>		real	<b>96.4±0.0</b>	<b>96.0±0.0</b>

## Robust vs occlusion and arbitrary orientation



## Wider gap in more challenging datasets

Method	Train data	36-char word accuracy per dataset			
		ArT 35,149	COCO 9,825	Uber 80,551	Total 125,525
ViTSTR-S	real	81.1±0.1	74.1±0.4	78.2±0.1	78.7±0.1
TRBA	real	82.5±0.2	77.5±0.2	81.2±0.3	81.3±0.2
ABINet	real	81.2±0.1	76.4±0.1	71.5±0.7	74.6±0.4
PARSeq <sub>N</sub>	real	83.0±0.2	77.0±0.2	82.4±0.3	82.1±0.2
PARSeq <sub>A</sub>	real	<b>84.5±0.1</b>	<b>79.8±0.1</b>	<b>84.5±0.1</b>	<b>84.1±0.0</b>