



2023

WHEN AND WHY TEST GENERATORS FOR DEEP LEARNING PRODUCE INVALID INPUTS: AN EMPIRICAL STUDY



**VINCENZO
RICCIO**

 @p1ndsvin



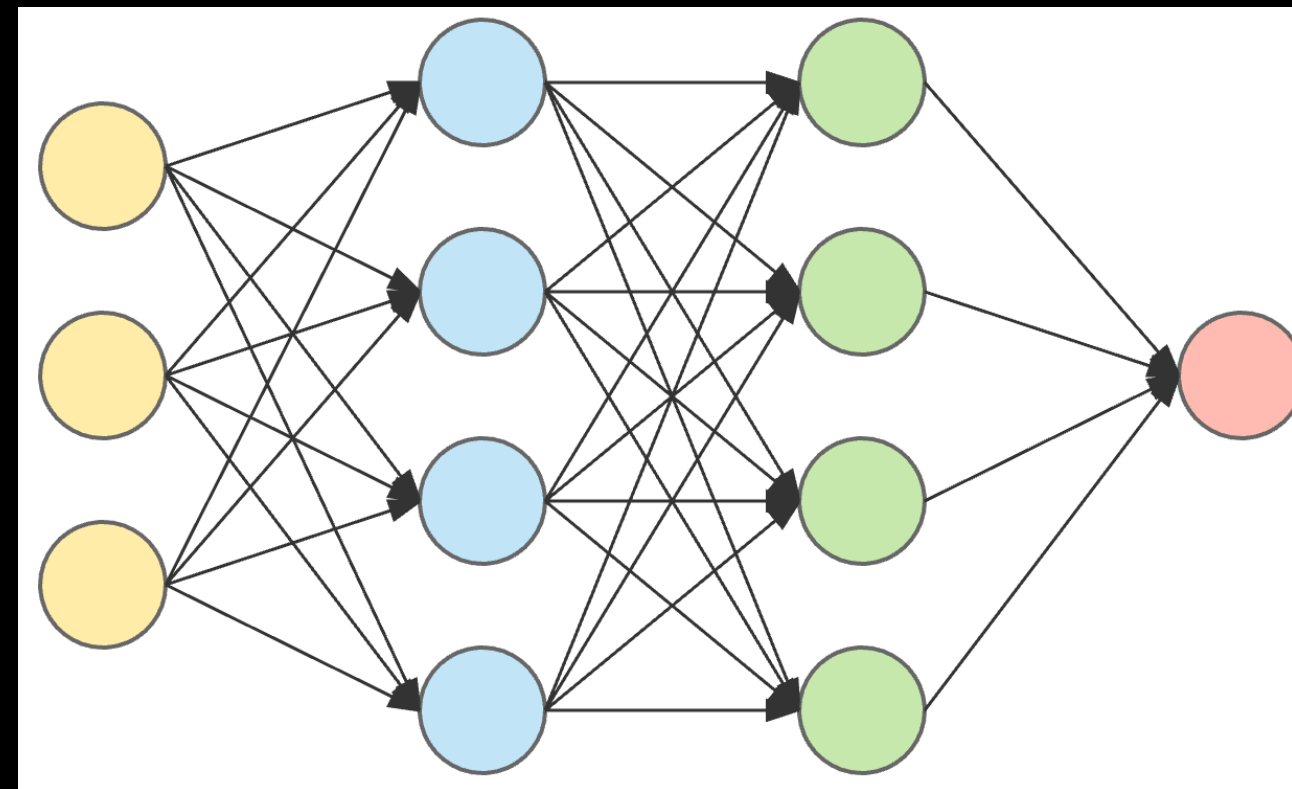
**PAOLO
TONELLA**

 @paolo_tonella

TRADITIONAL DL MODEL ASSESSMENT



**ORIGINAL
DATASET**



**DL MODEL
UNDER TEST**

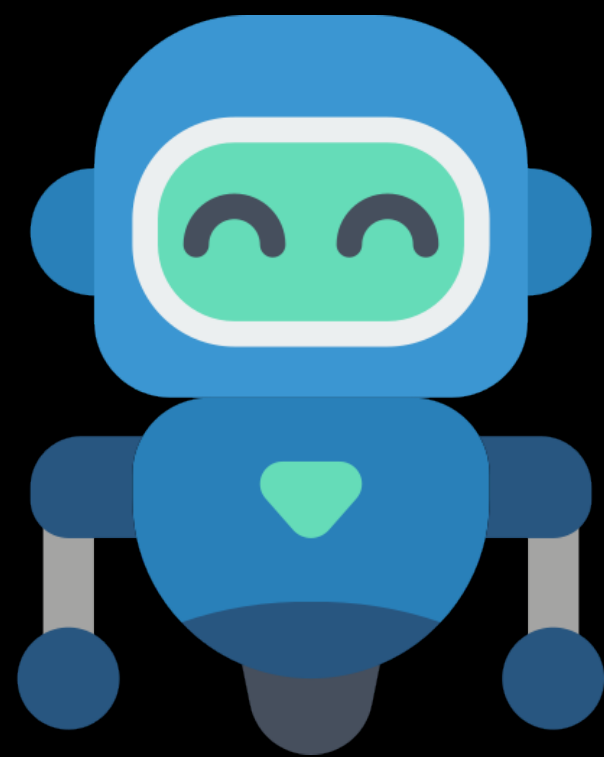


**ACC = 0.95
MSE = 0.01**

**PERFORMANCE
METRIC**

What is the performance of a DL model for inputs beyond its original dataset?

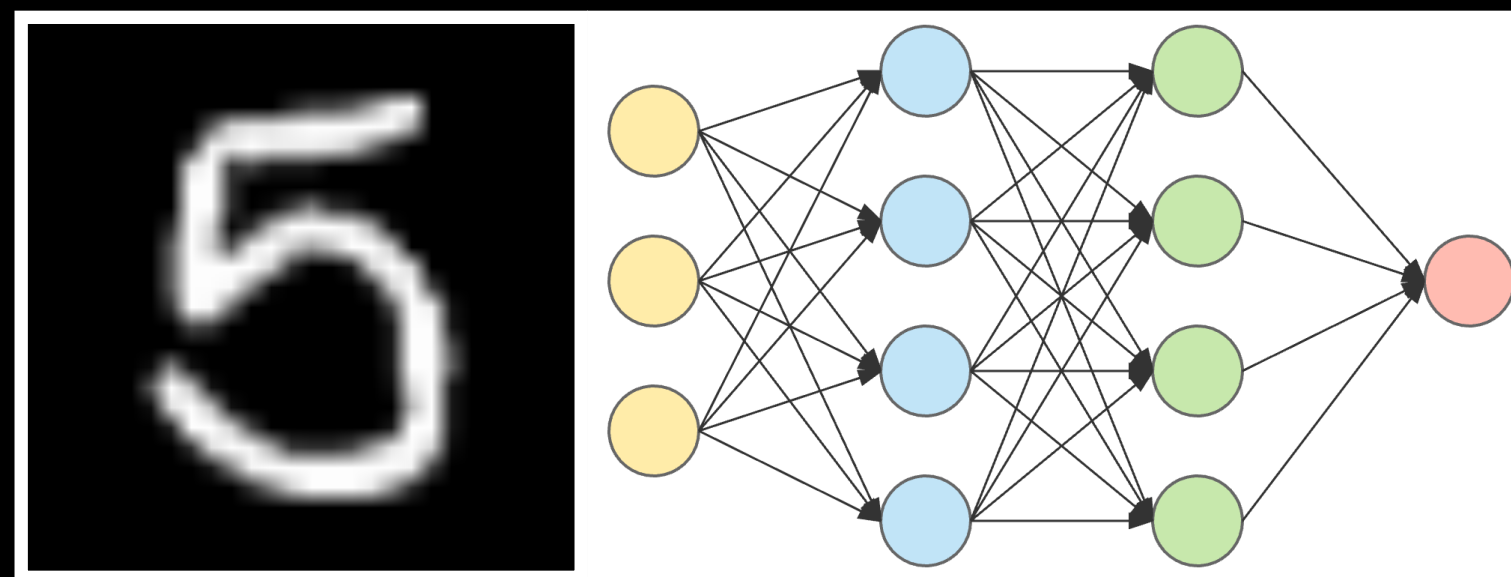
AUTOMATED TEST INPUT GENERATION FOR DL MODELS



**TEST
GENERATOR**

Target Label

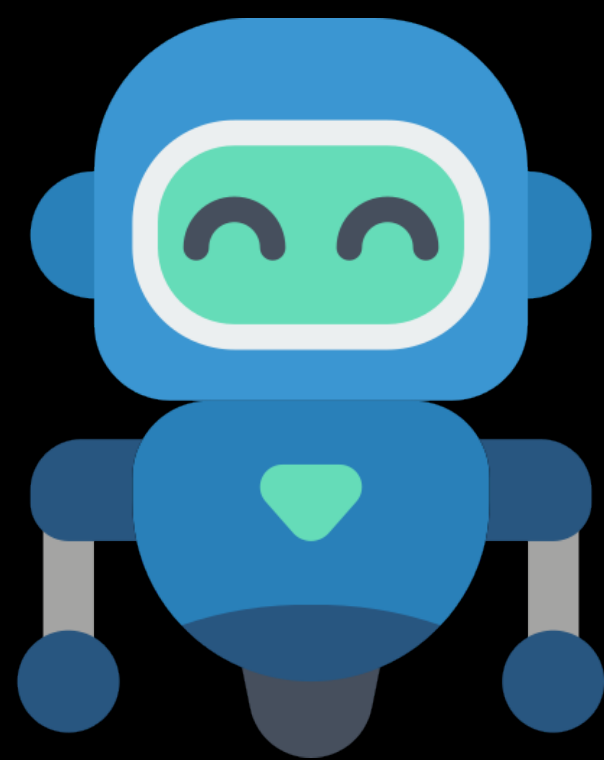
5



Predicted Label

5 ✓

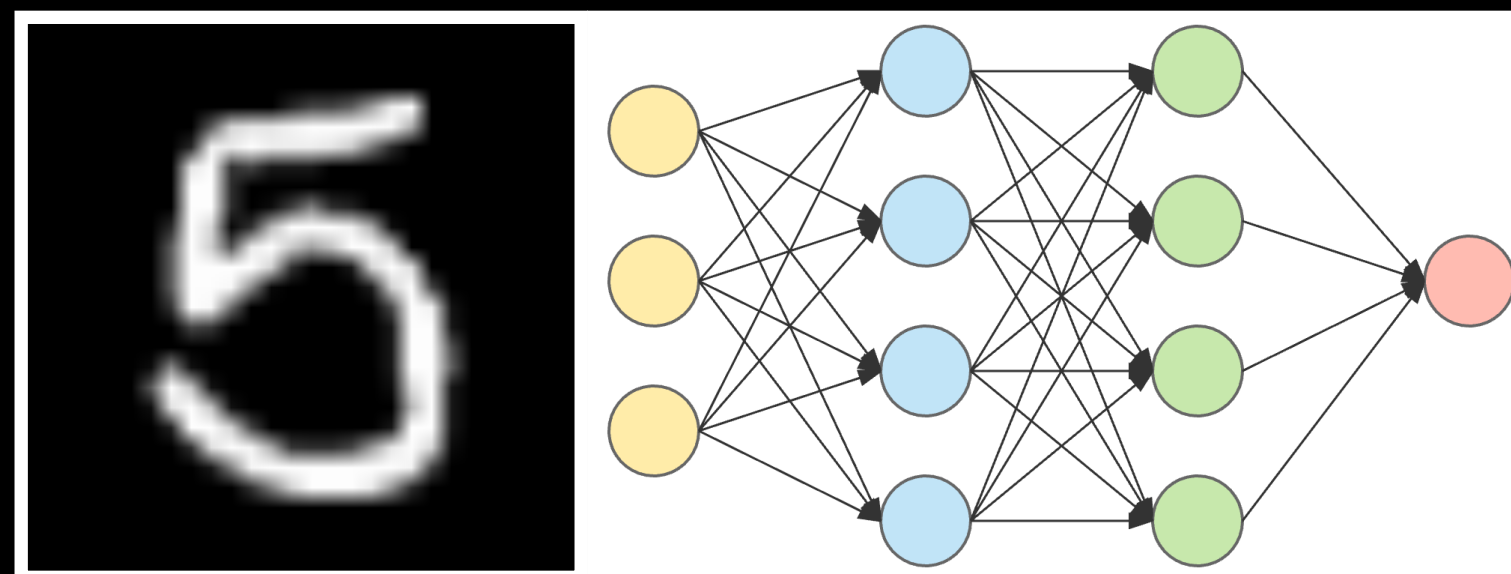
AUTOMATED TEST INPUT GENERATION FOR DL MODELS



**TEST
GENERATOR**

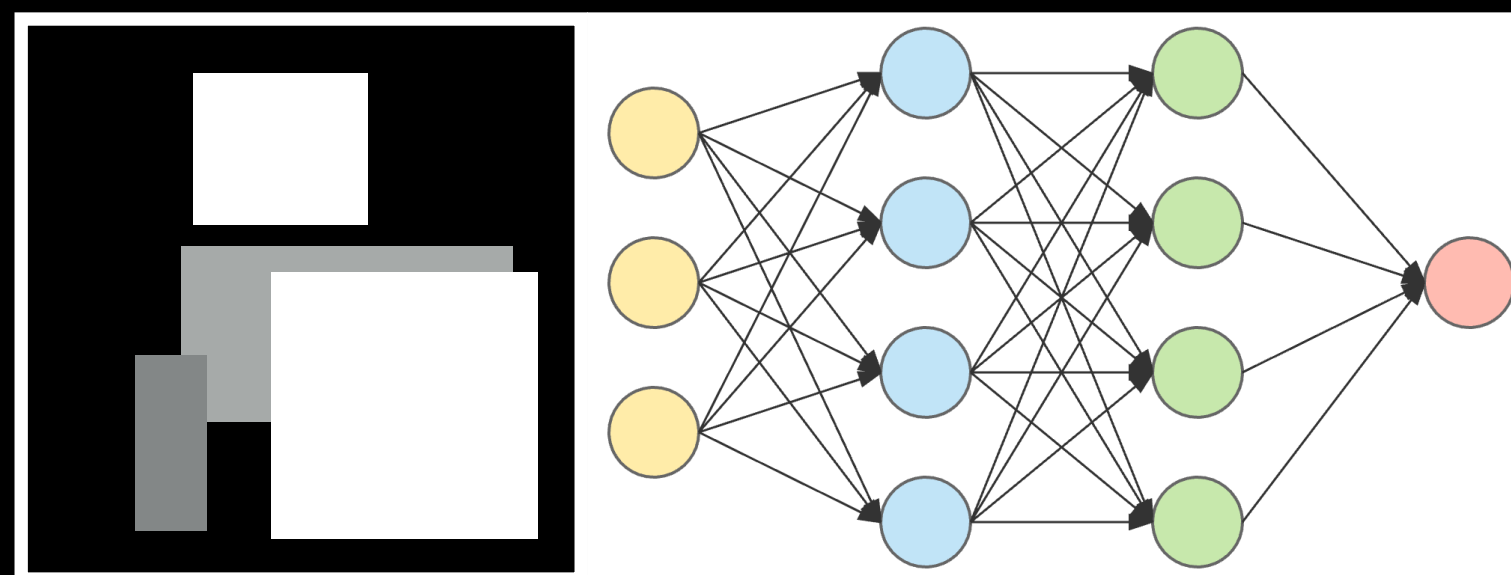
Target Label

5



Predicted Label

5 ✓



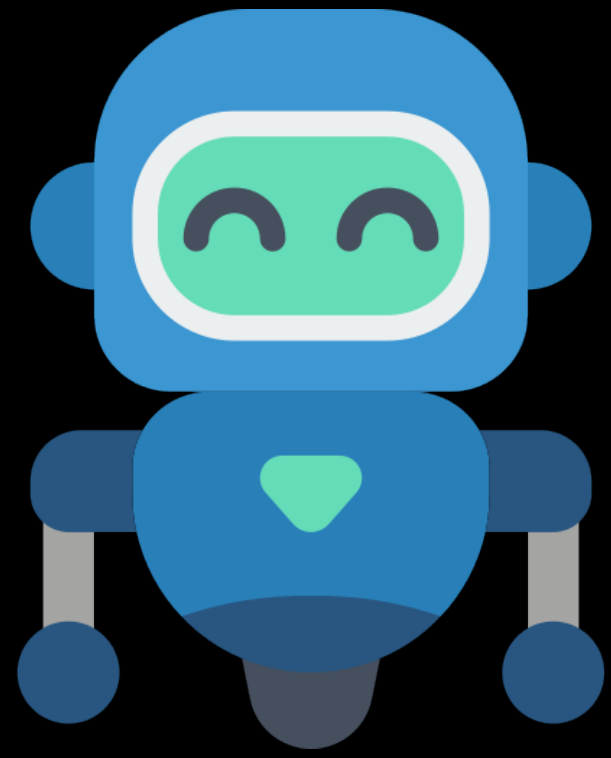
Predicted Label

6 ✗

Problem #1:

invalid inputs, not
recognisable by domain
experts in the input domain

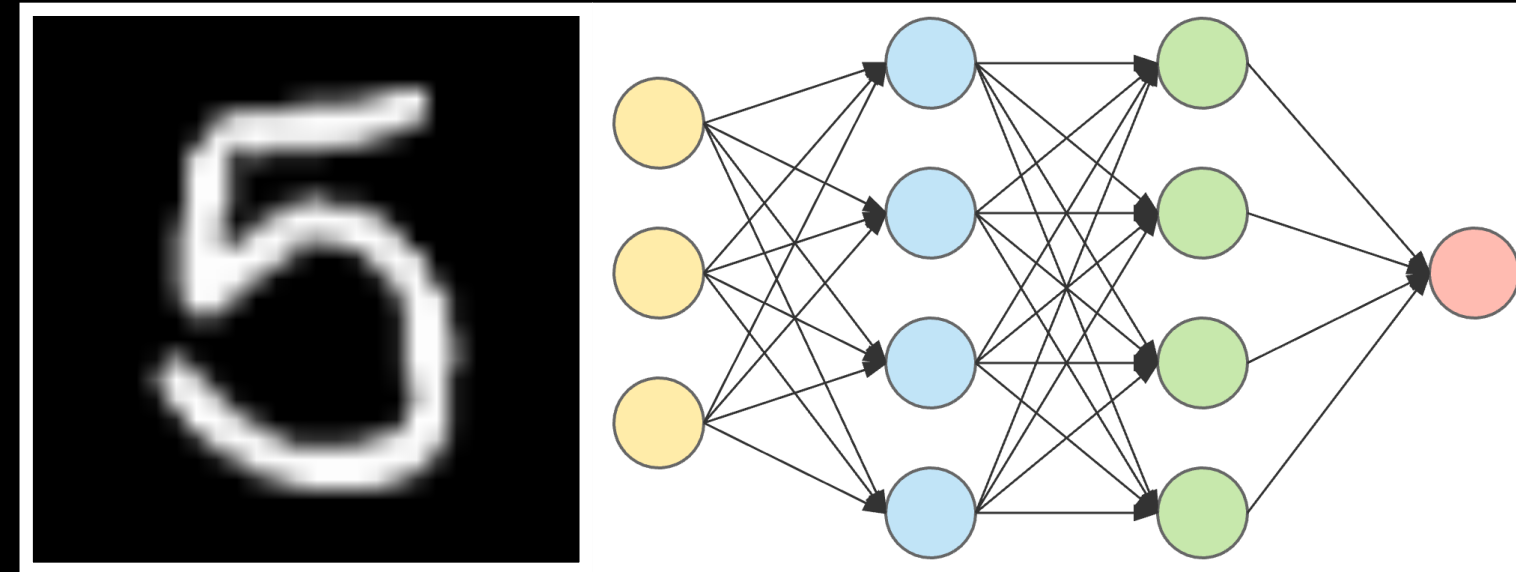
AUTOMATED TEST INPUT GENERATION FOR DL MODELS



**TEST
GENERATOR**

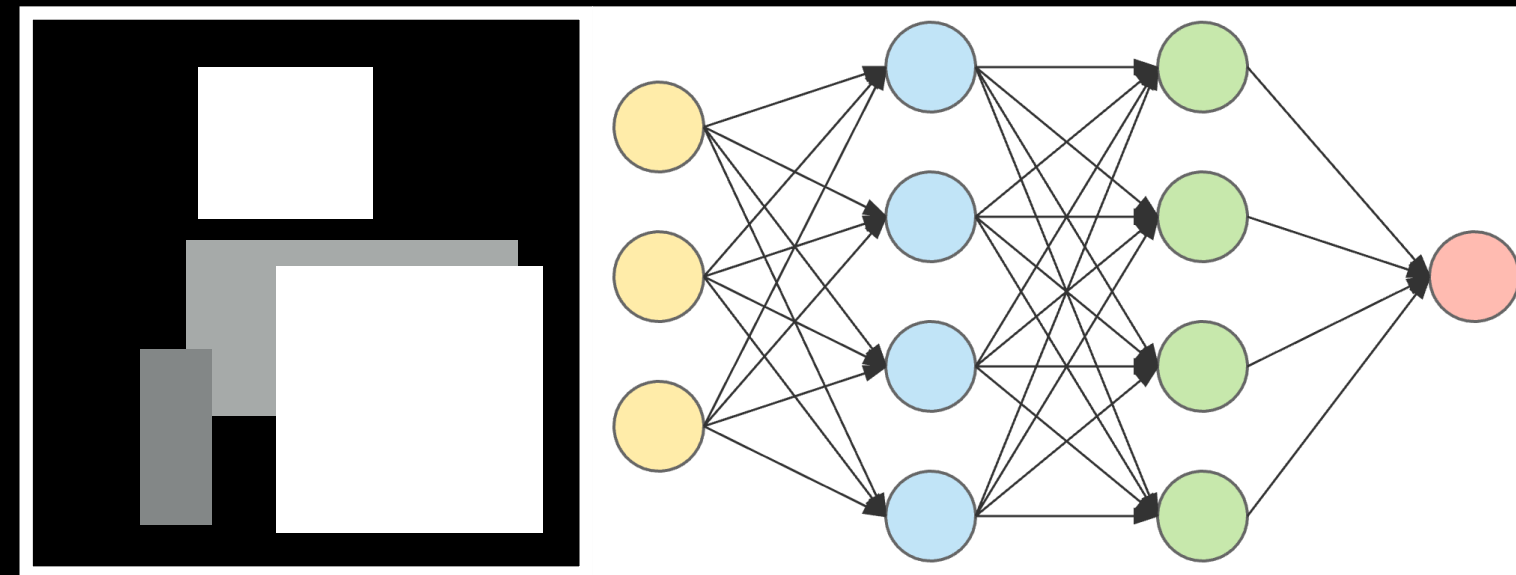
Target Label

5



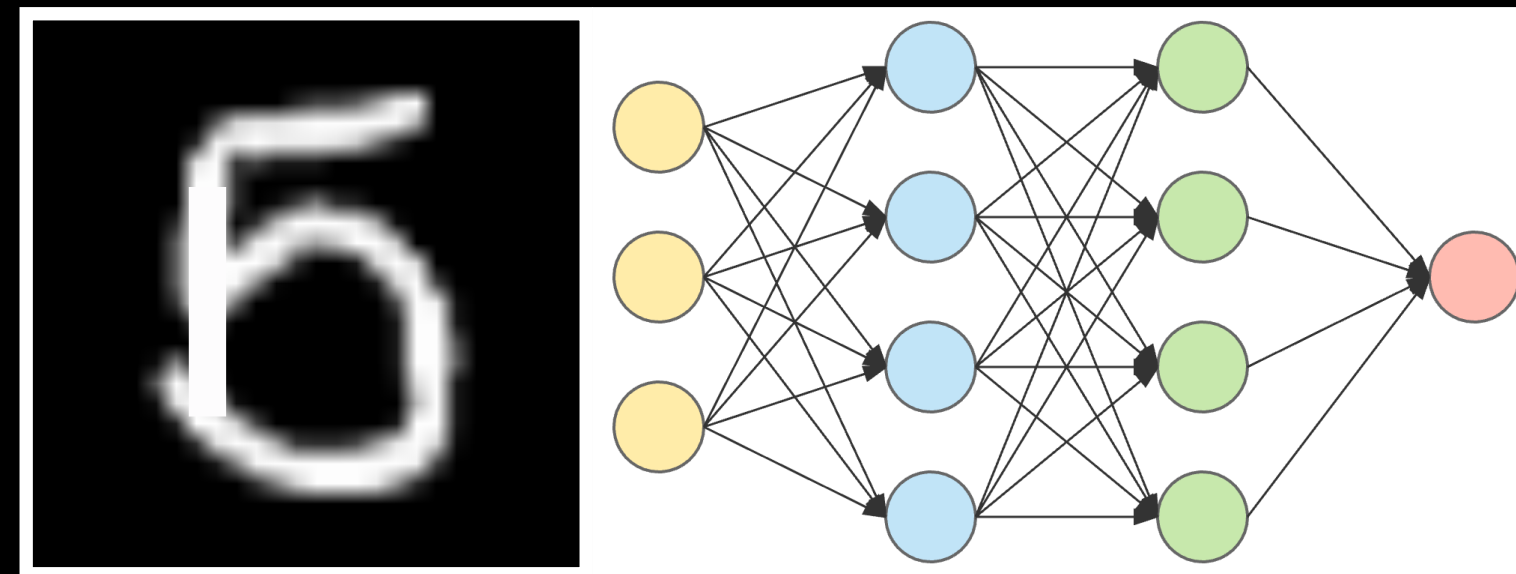
Predicted Label

5 ✓



Predicted Label

6 ✗



Predicted Label

6 ✗

Problem #1:

invalid inputs, not
recognisable by domain
experts in the input domain

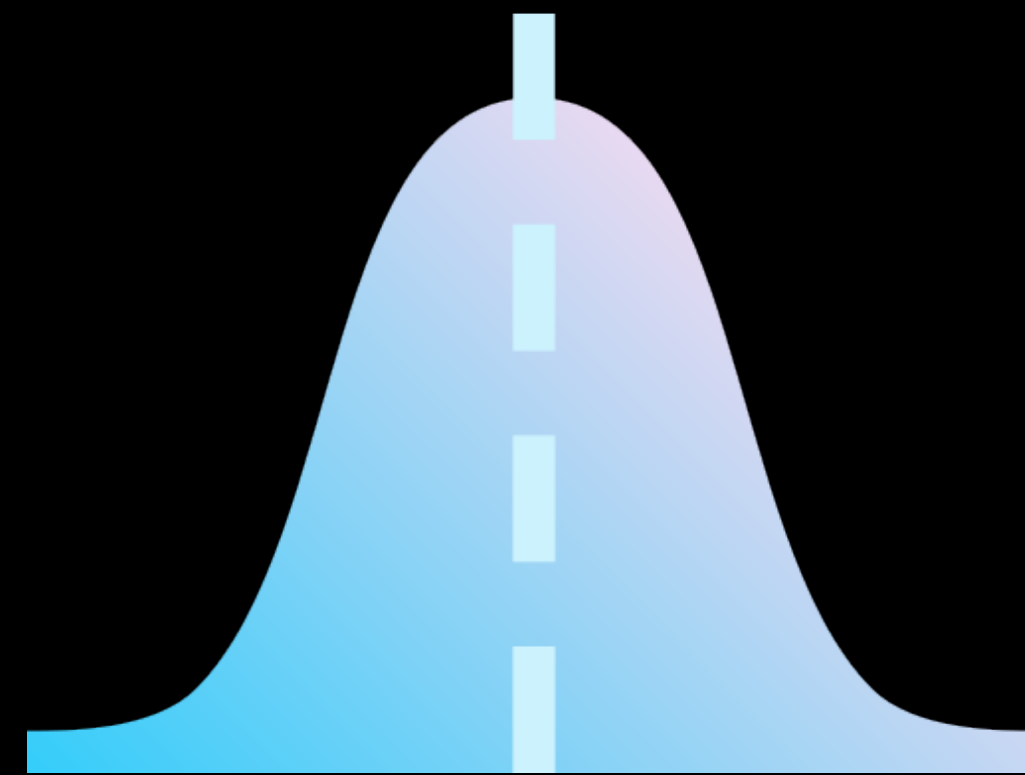
Problem #2:

original label is not
preserved

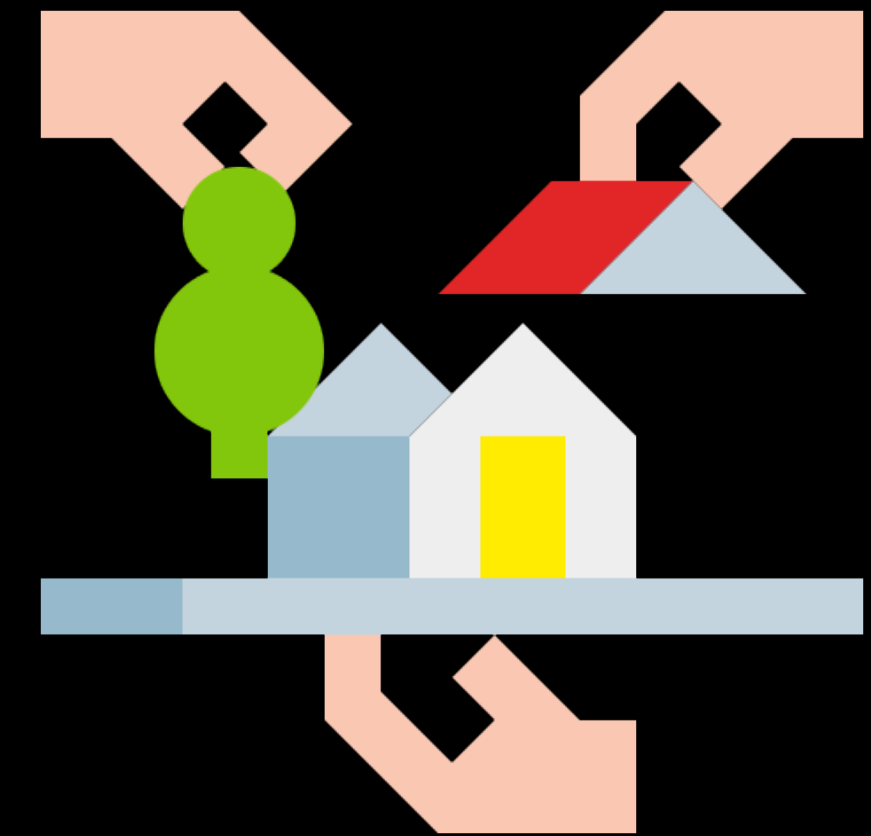
TEST INPUT GENERATION APPROACHES



**RAW INPUT
MANIPULATION**



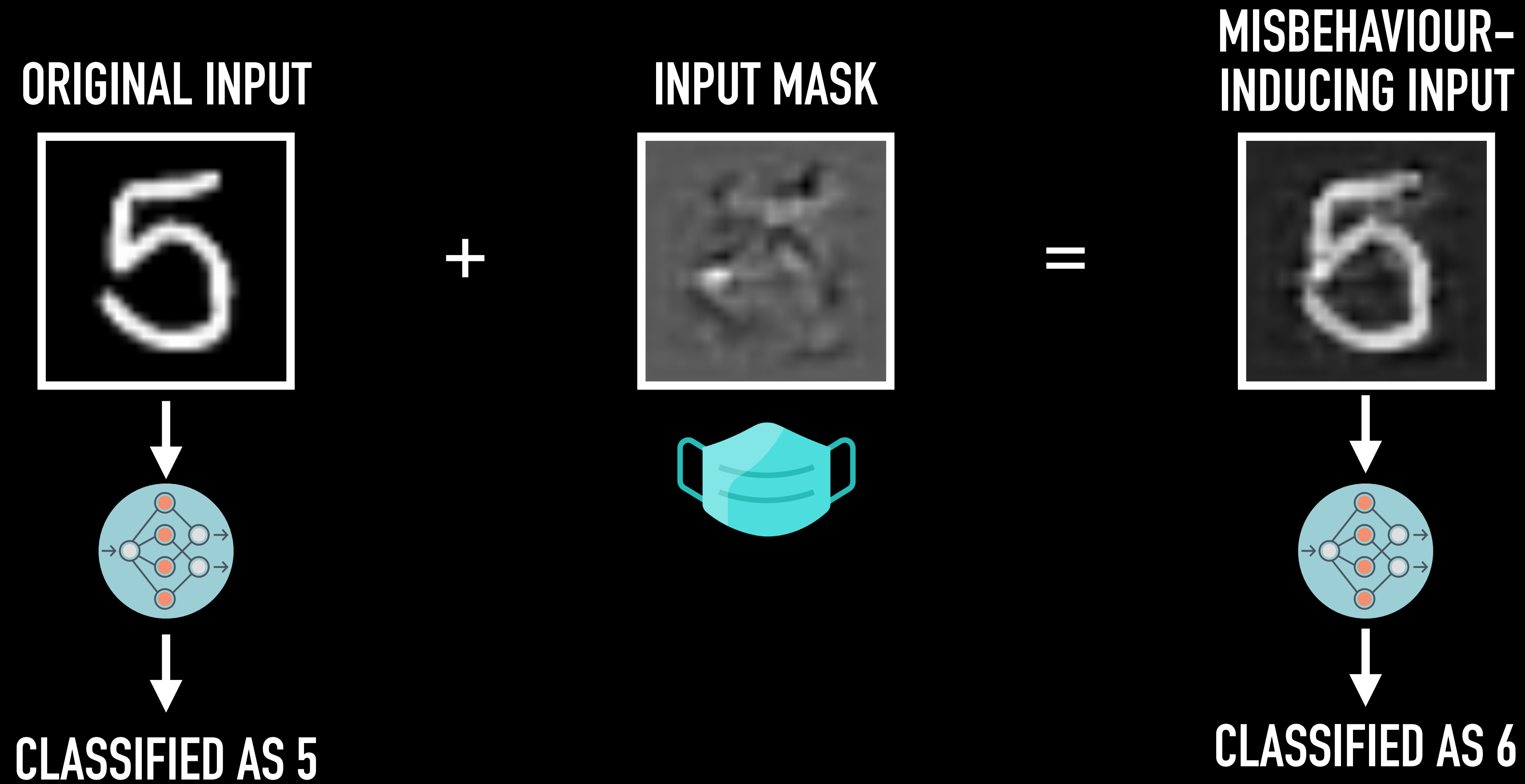
**GENERATIVE DL
MODELS**



**MODEL-BASED INPUT
MANIPULATION**

DLFUZZ

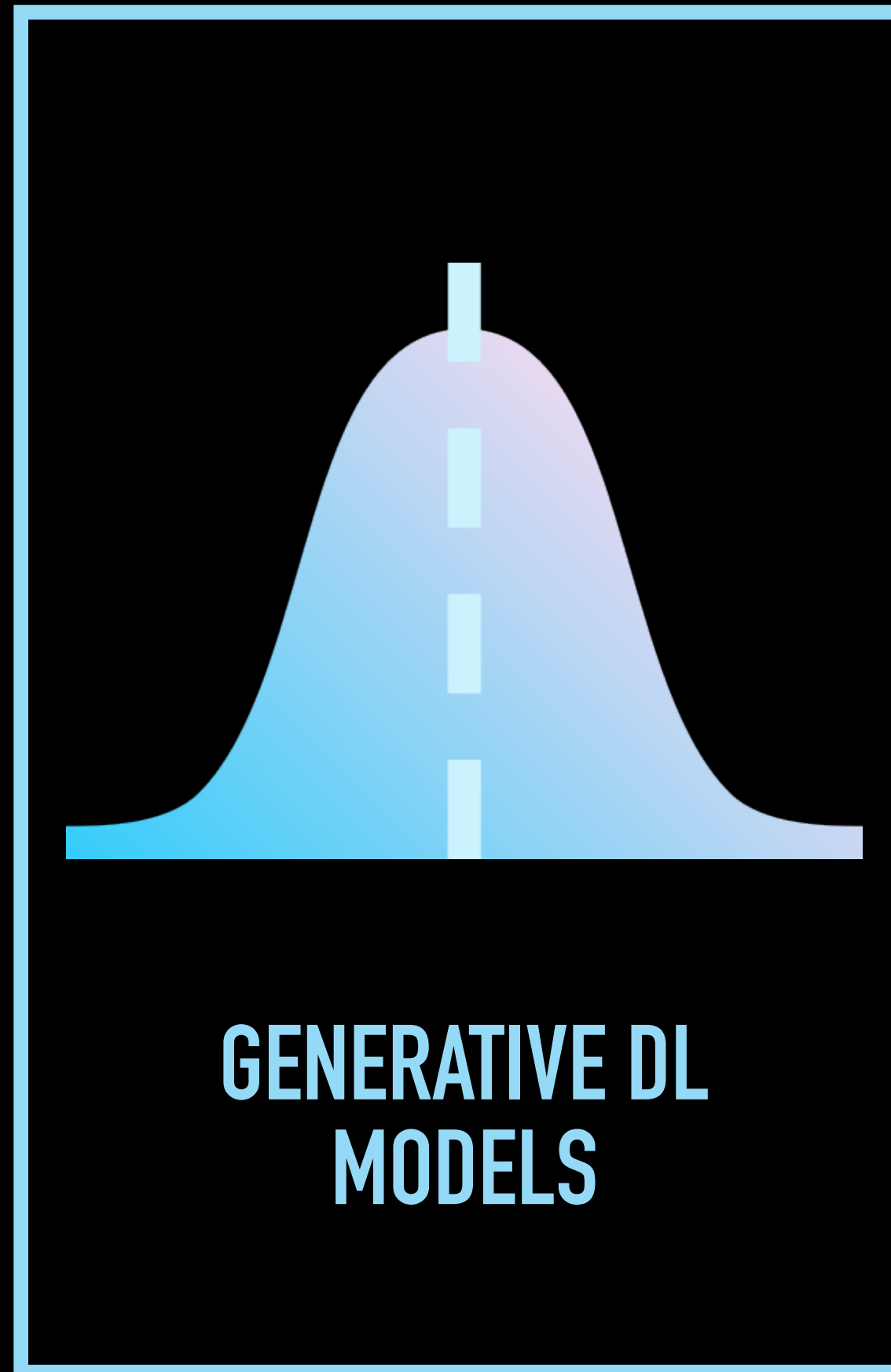
Guo et al., FSE 2018



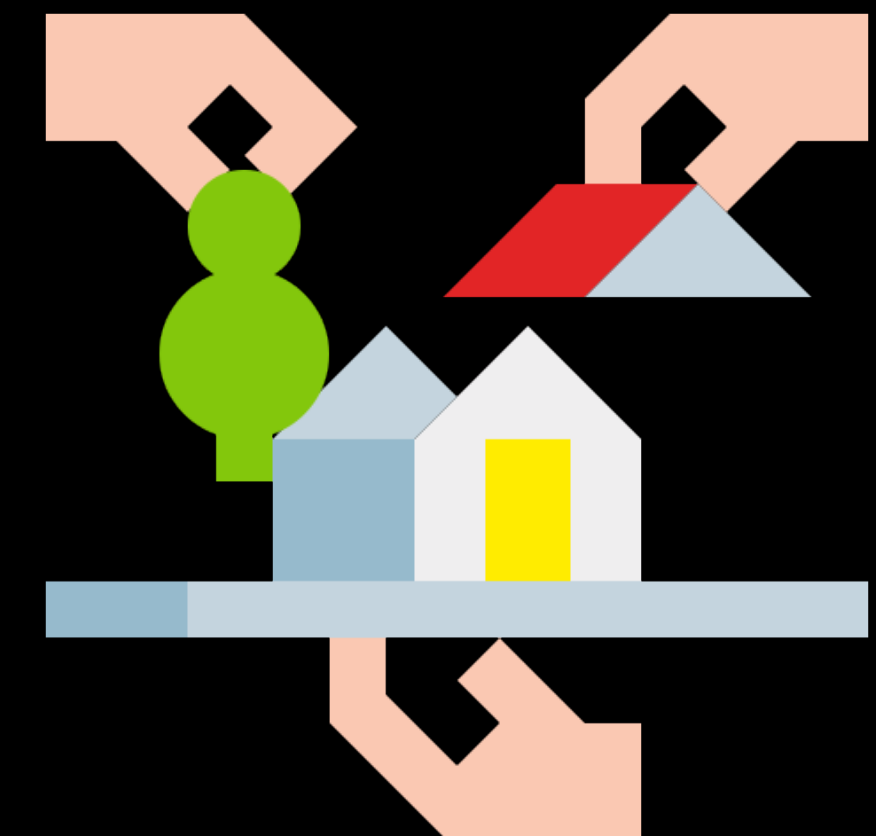
TEST INPUT GENERATION APPROACHES



**RAW INPUT
MANIPULATION**

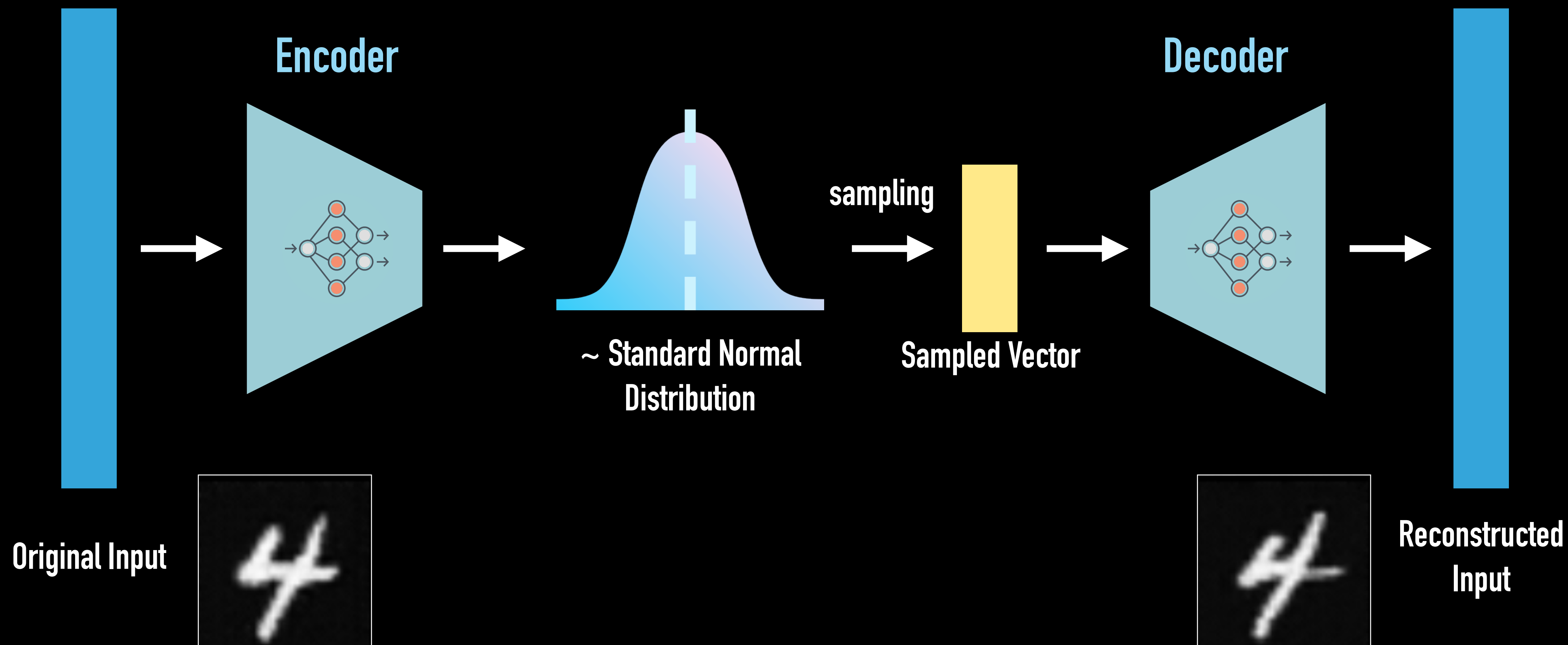


**GENERATIVE DL
MODELS**

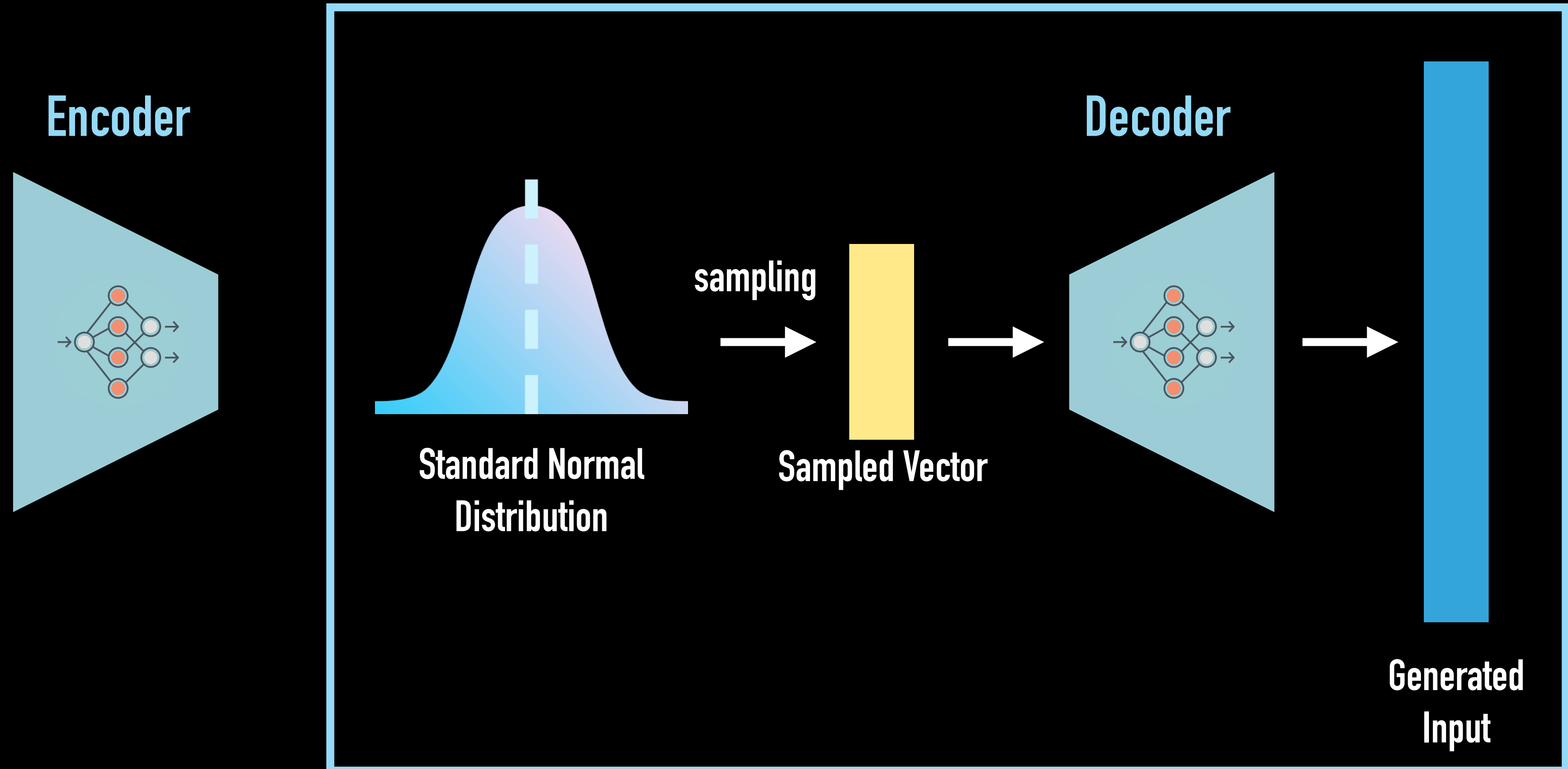


**MODEL-BASED INPUT
MANIPULATION**

VARIATIONAL AUTOENCODER (VAE)

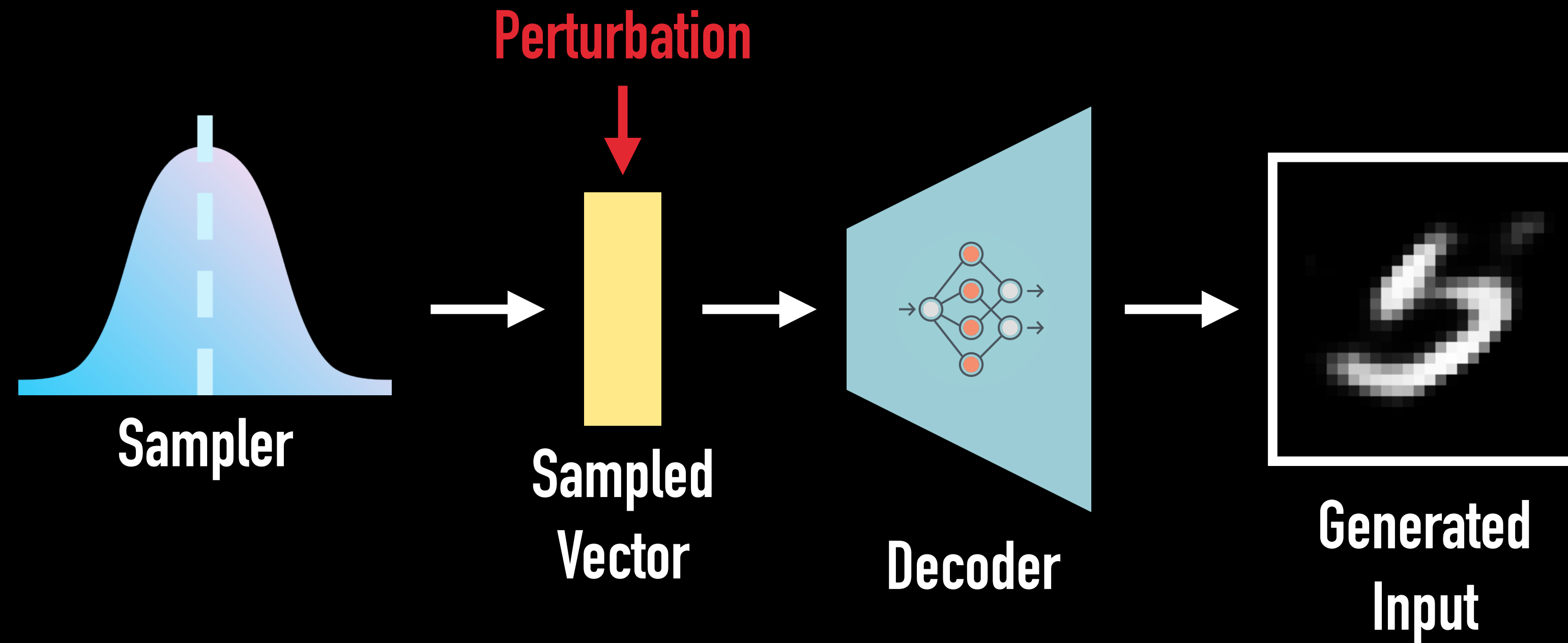


VARIATIONAL AUTOENCODER (VAE)



SINVAD

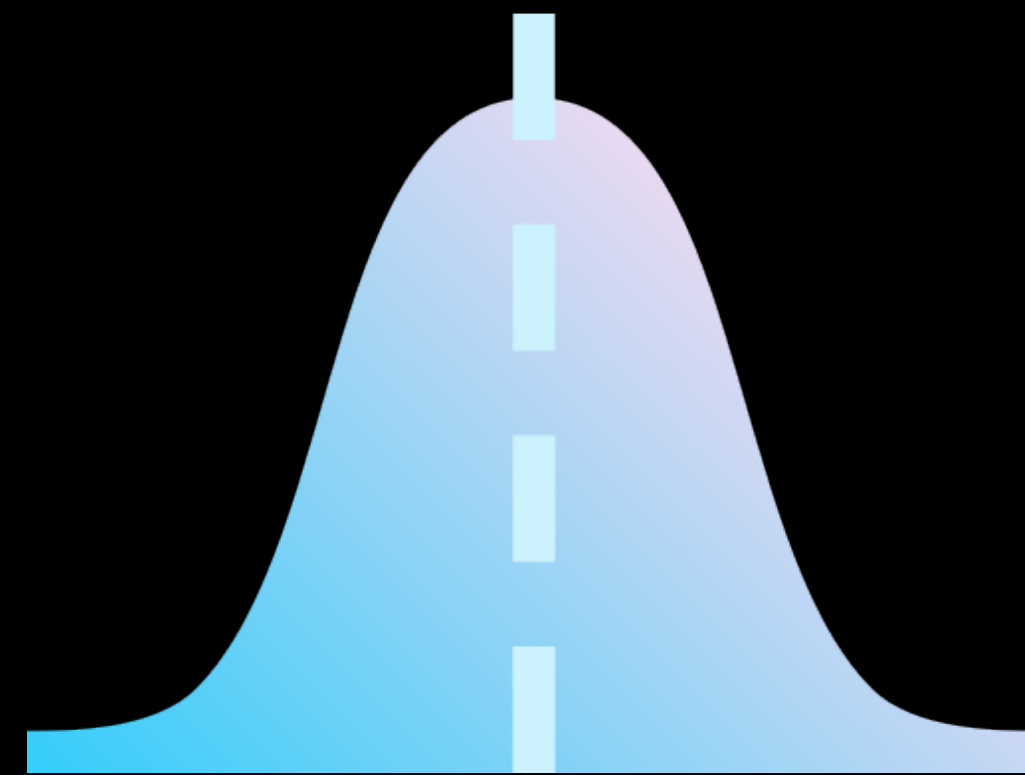
Kang et al., ICSE 2018



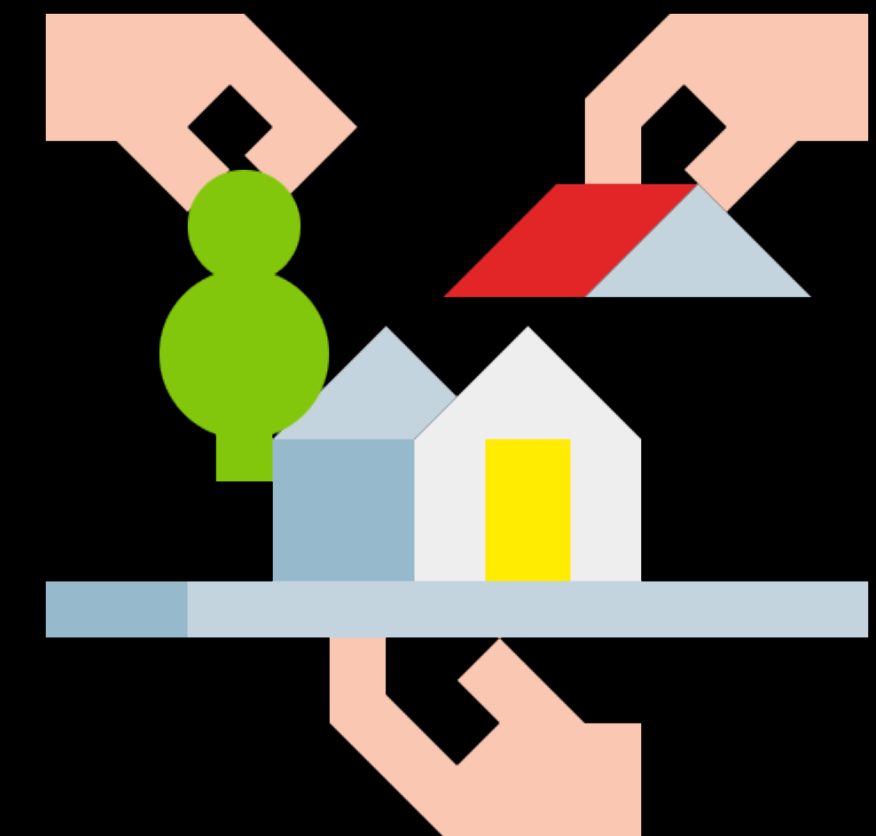
TEST INPUT GENERATION APPROACHES



**RAW INPUT
MANIPULATION**



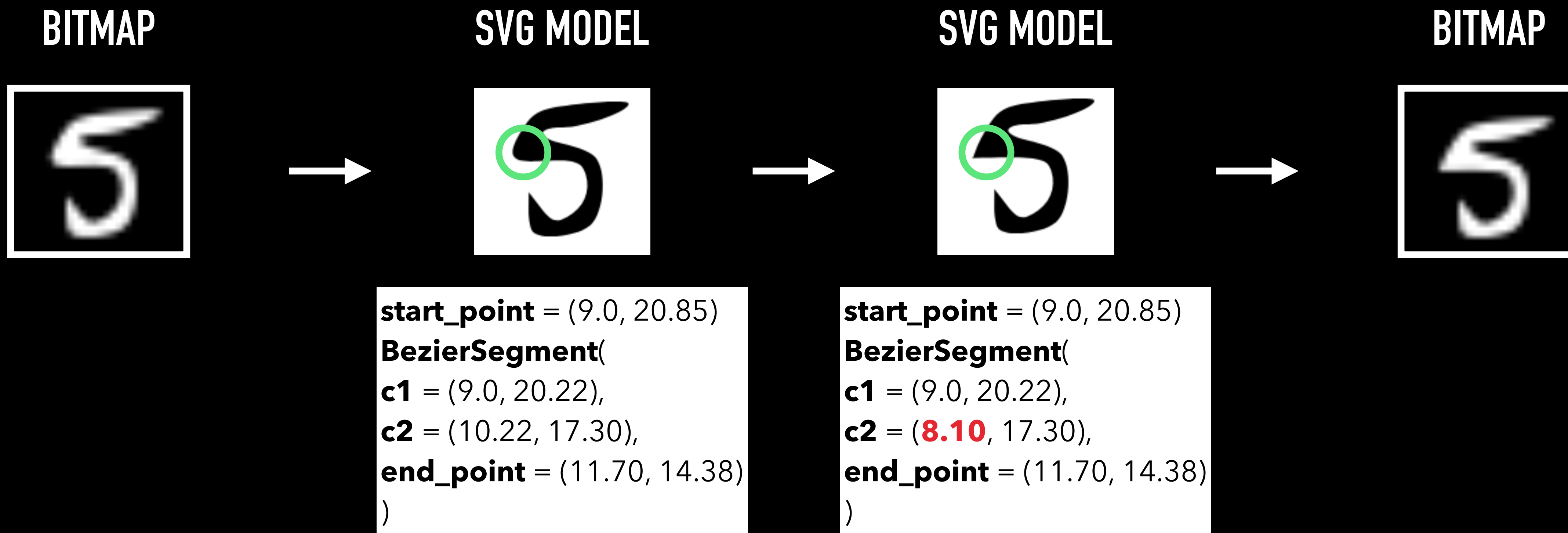
**GENERATIVE DL
MODELS**



**MODEL-BASED INPUT
MANIPULATION**

DEEPJANUS

Riccio and Tonella, FSE 2020



**ARE INPUTS PRODUCED
BY TEST INPUT
GENERATORS VALID?**

**ARE THE GENERATED
VALID INPUTS
LABEL-PRESERVING?**

HUMAN ASSESSMENT

220 HUMAN ASSESSORS FROM

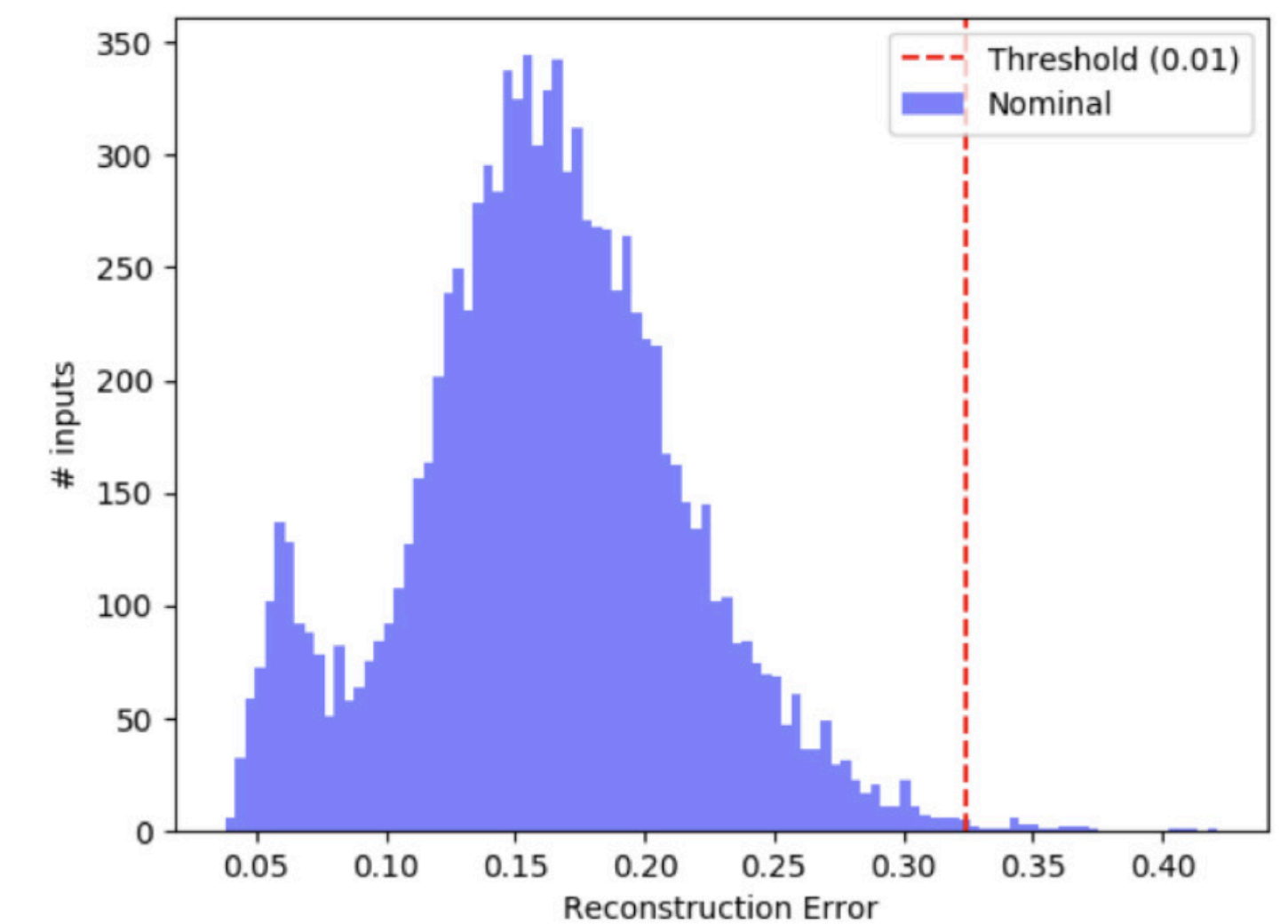
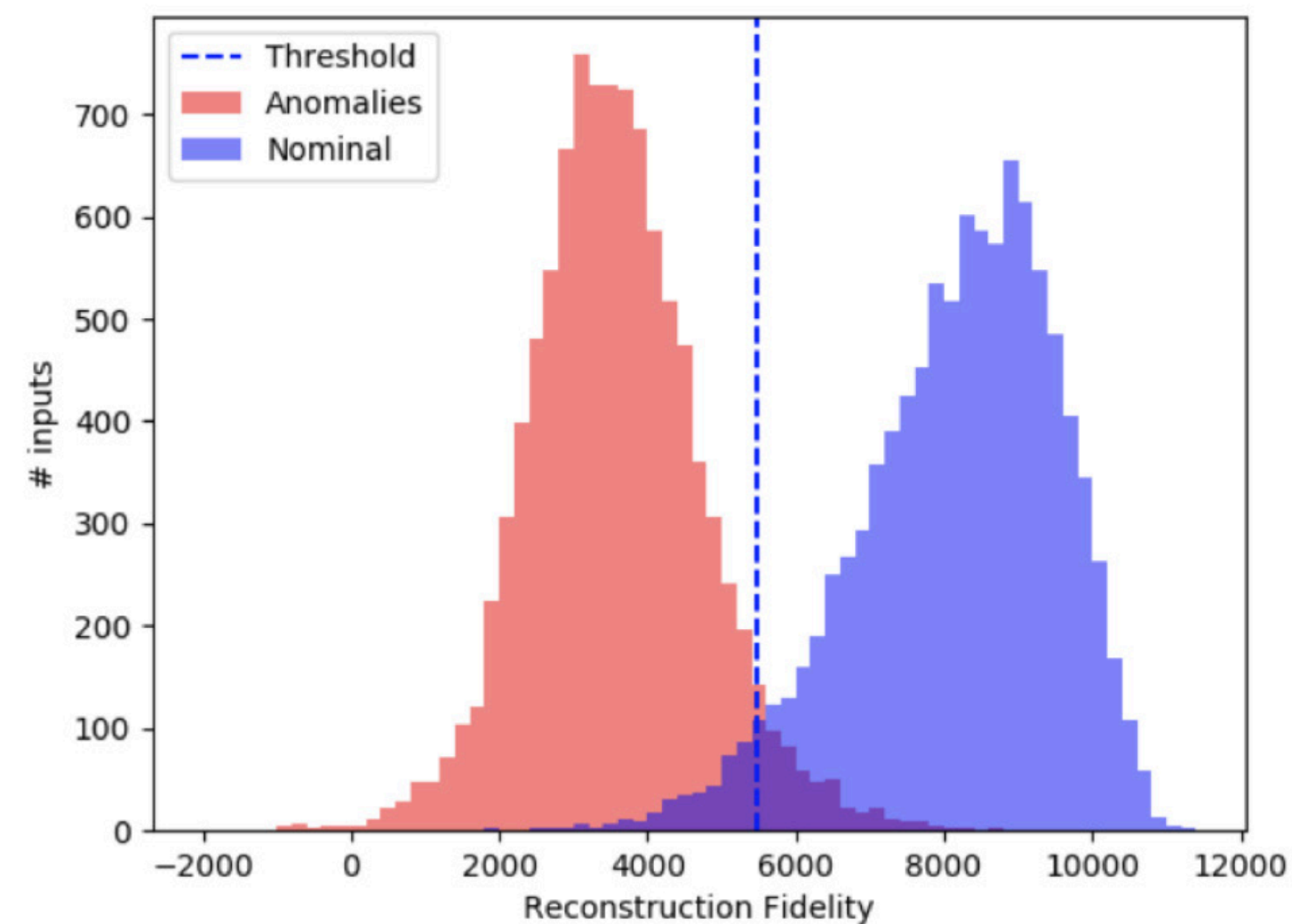
amazon mechanical turk



AUTOMATED ASSESSMENT

DAIV
[DOLA ET AL., ICSE 2021]

SELFORACLE
[STOCCO ET AL., ICSE 2020]

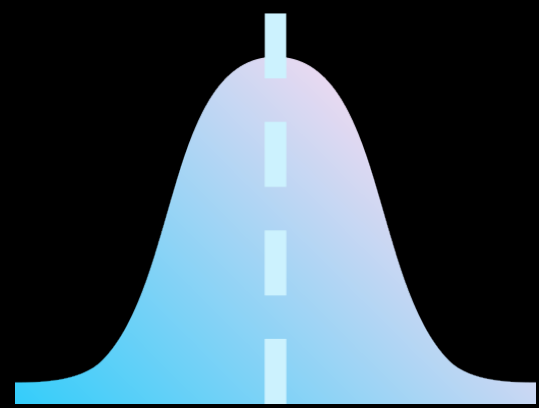


EMPIRICAL STUDY: TEST GENERATORS



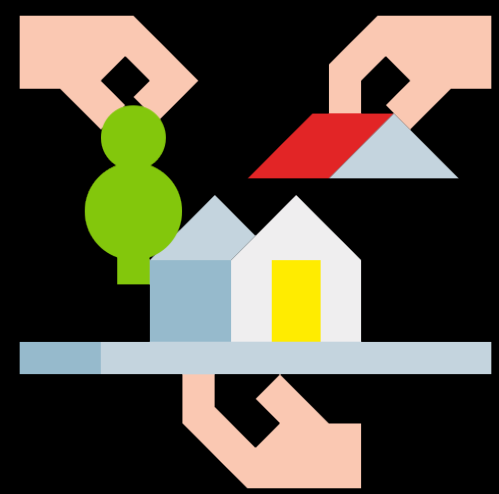
RAW INPUT MANIPULATION

- ▶ DEEXPLORE [PEI ET AL., SOSP 2017]
- ▶ DLFUZZ [GUO ET AL., FSE 2018]



GENERATIVE DL MODELS

- ▶ SINVAD [KANG ET AL., ICSE 2018]
- ▶ FEATURE PERTURBATIONS [DUNN ET AL., ISSTA 2021]



MODEL-BASED INPUT MANIPULATION

- ▶ DEEPIANUS [RICCIO AND TONELLA, FSE 2020]

DATASETS



MNIST



SVHN



IMAGENET-1K

HUMAN ASSESSMENT

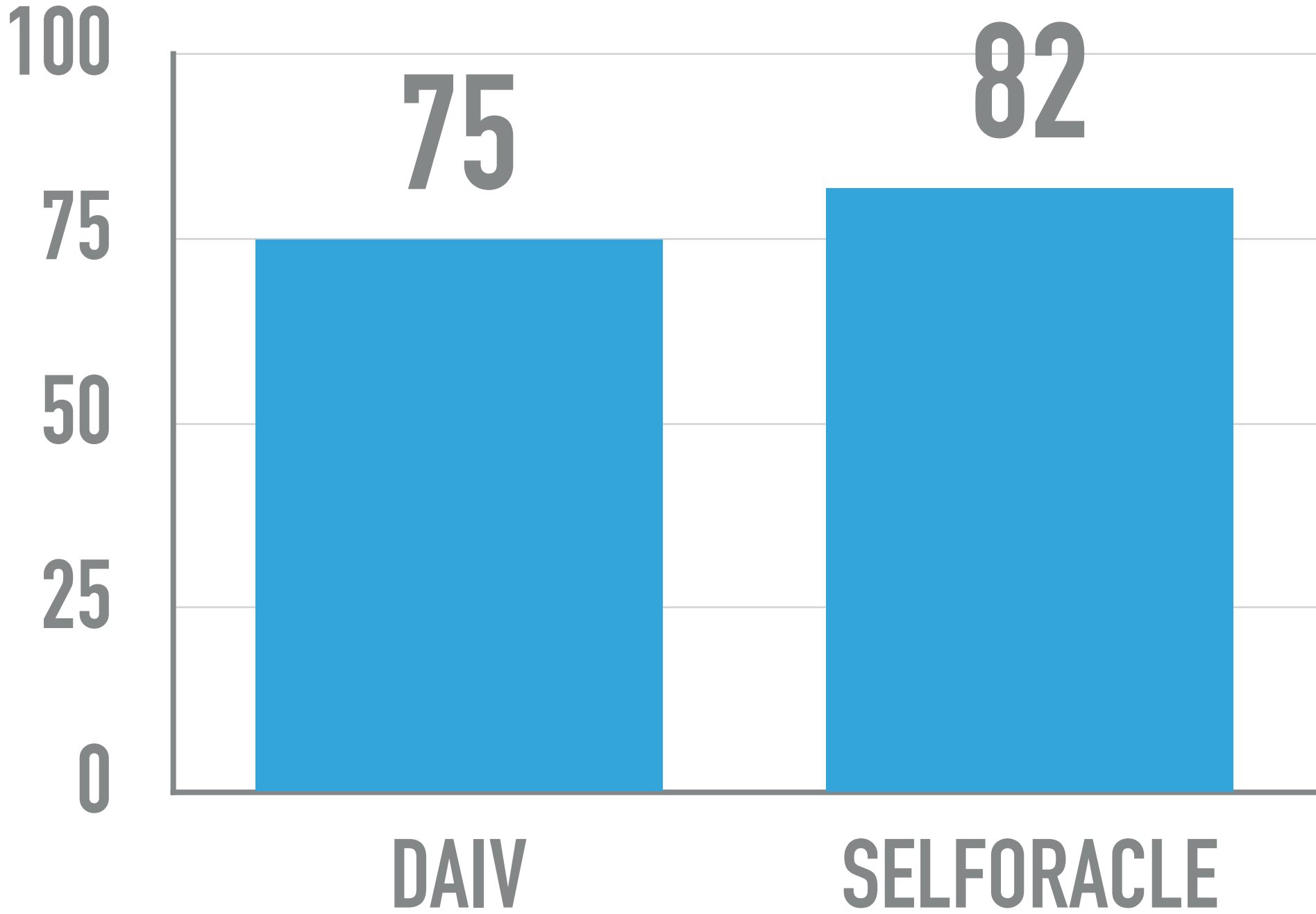


**88% VALID
INPUTS**

**69%
PRESERVED
LABELS**

AUTOMATED ASSESSMENT

**% AGREEMENT WITH HUMAN
ASSESSMENT**



LESSONS LEARNT



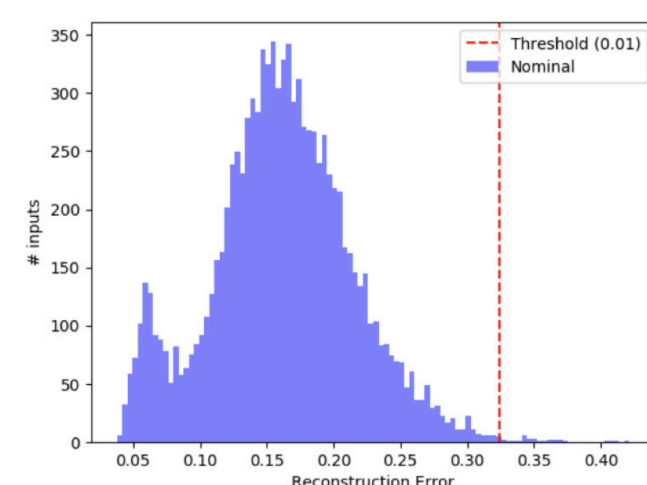
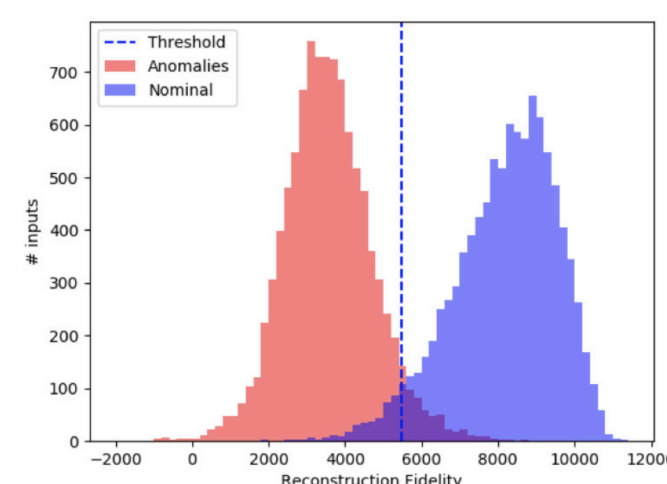
Too aggressive raw data manipulations lead to invalid inputs



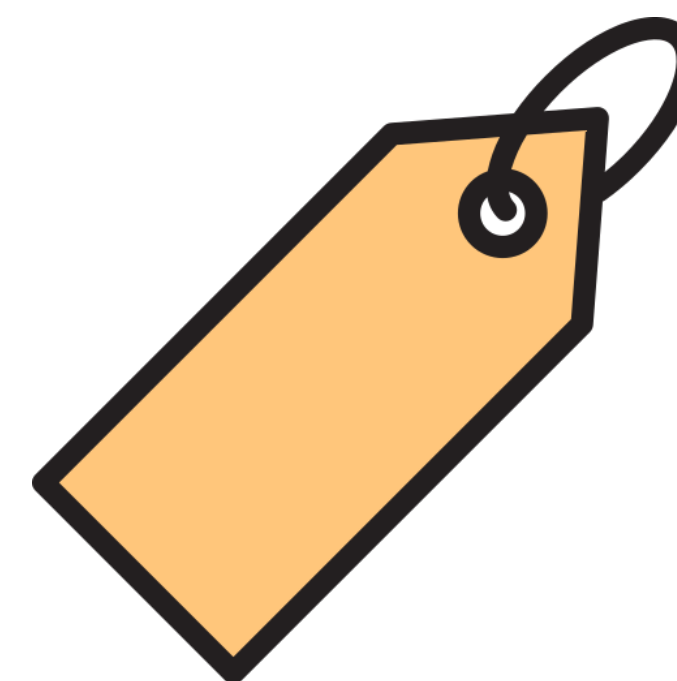
Generative DL models should carefully explore the latent space



Model-based techniques require high-quality model representations

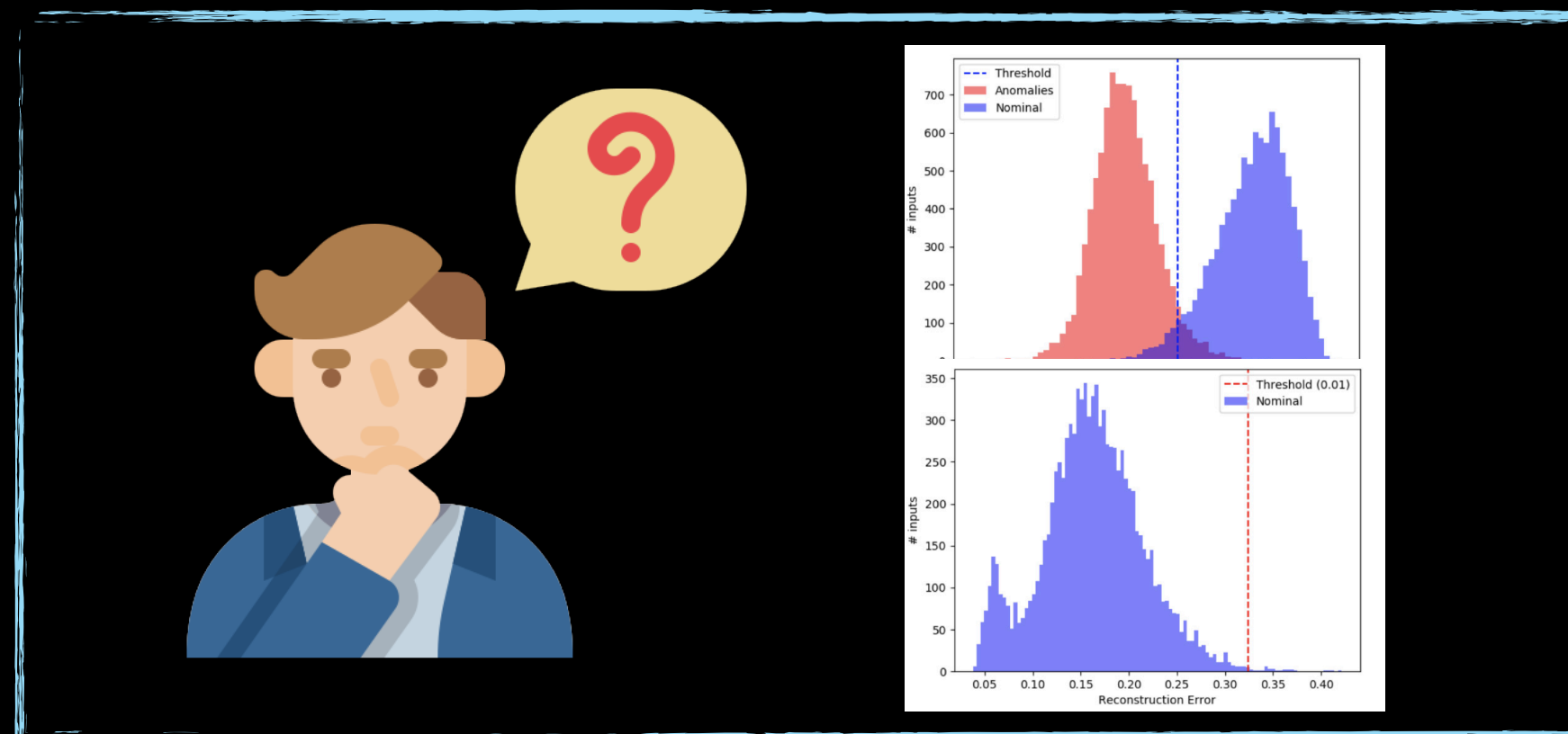
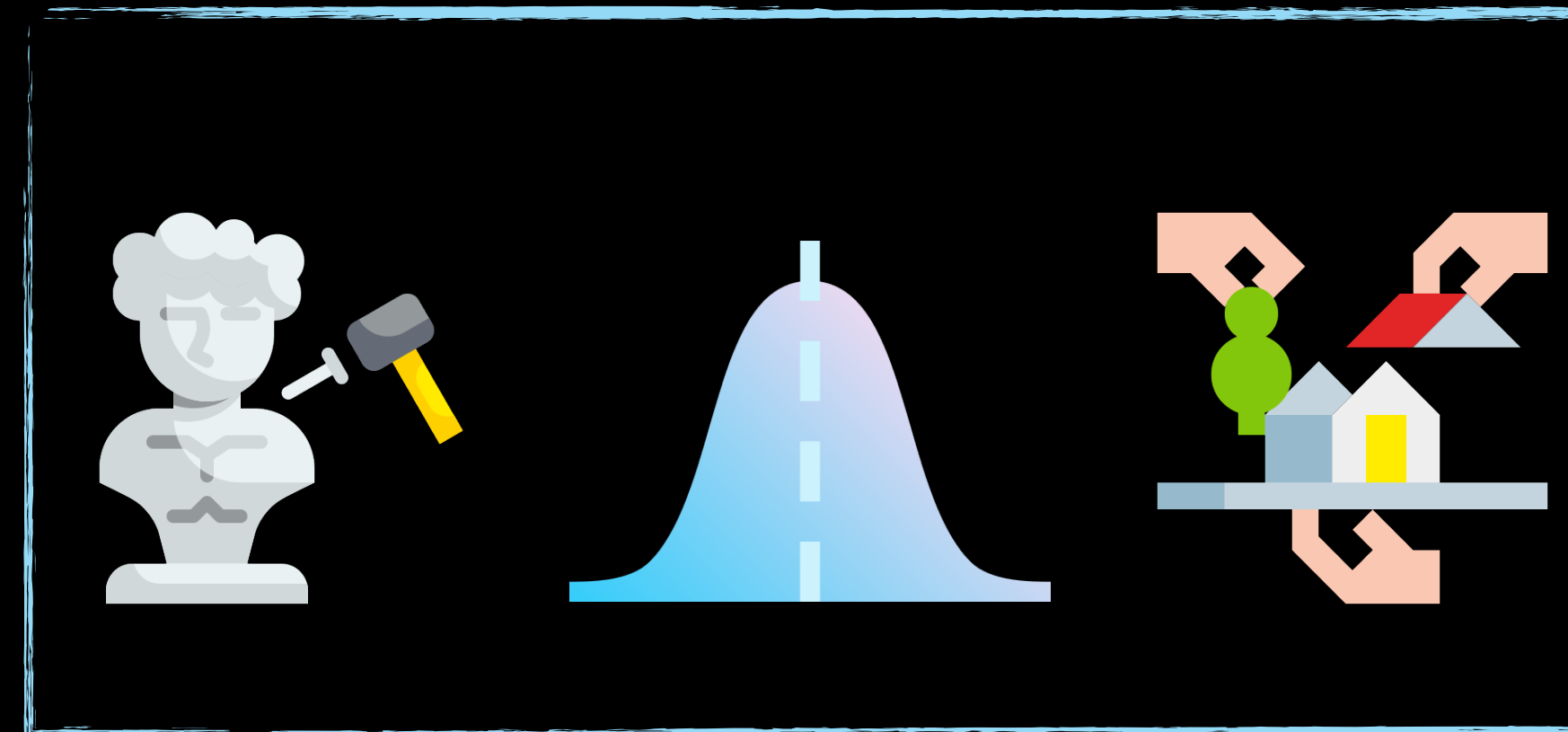
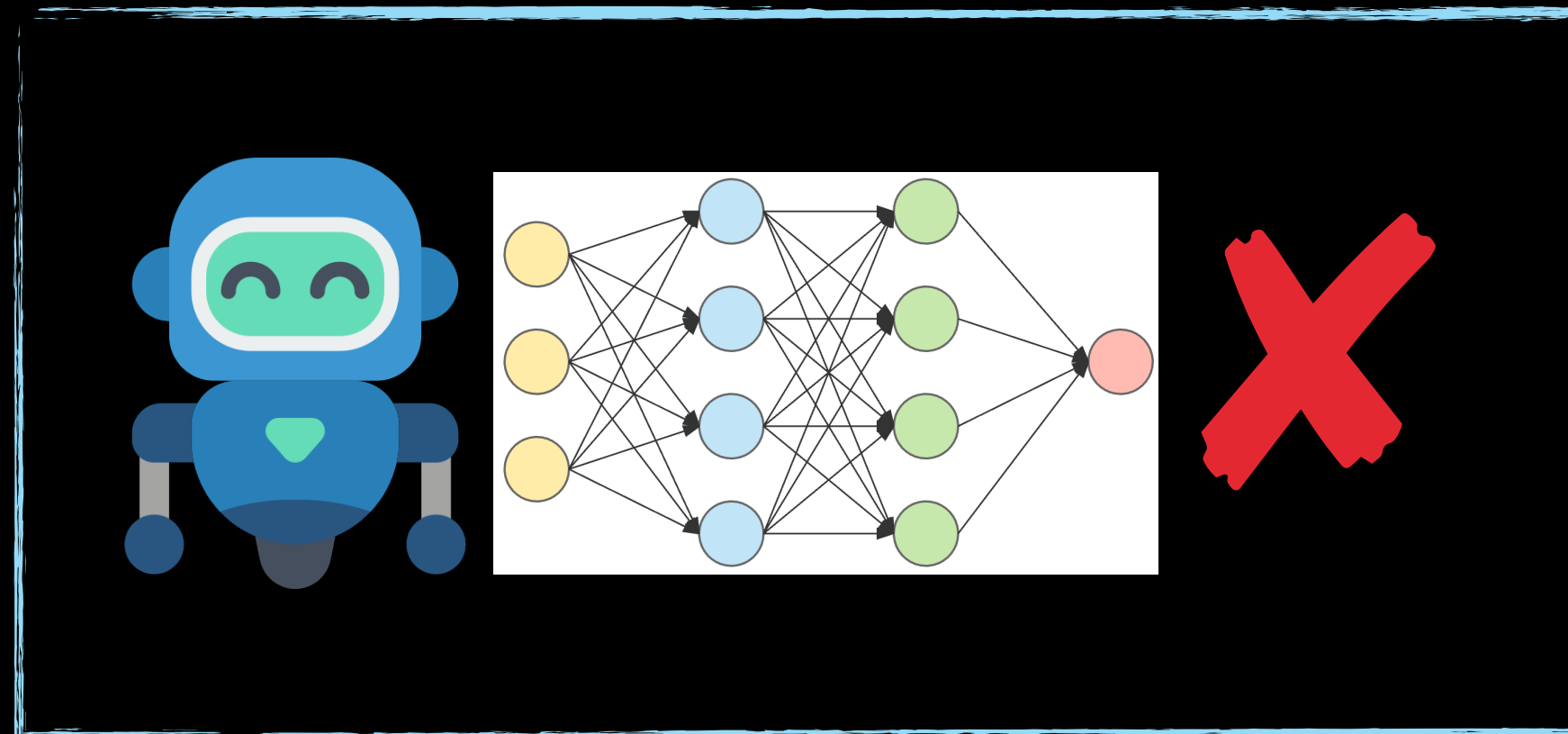


Automated validators check in-distribution



Label preservation is mostly overlooked

SUMMARY






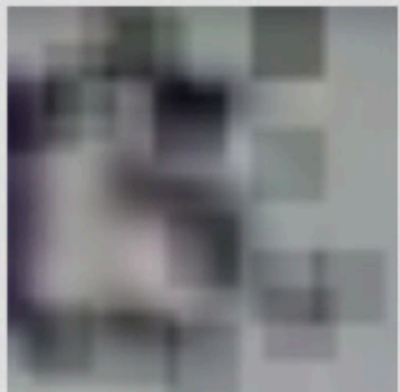
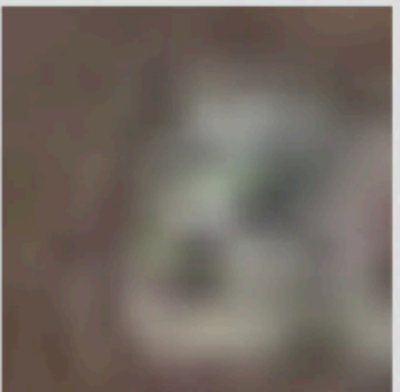
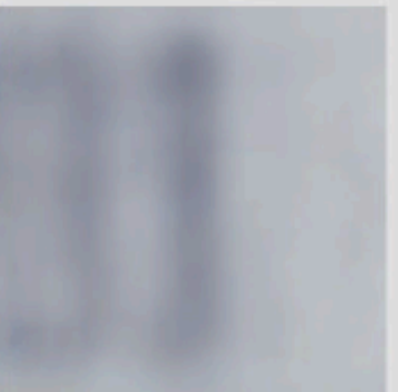
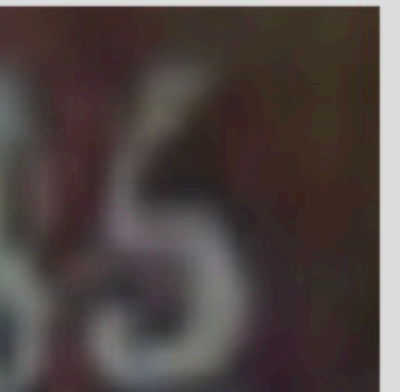
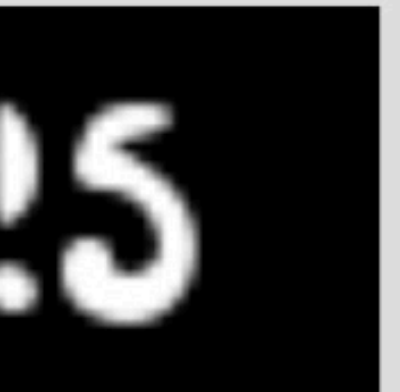


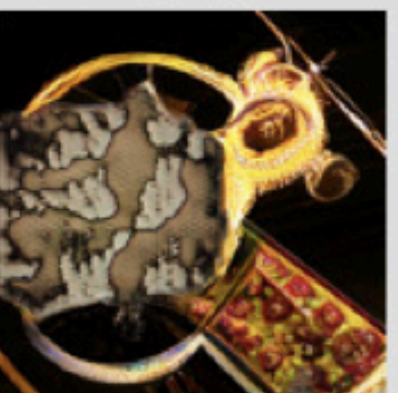




`testingautomated-usi/
tig-validity-icse23`

EXTRA SLIDES

Dataset	Tool	% Valid			% Hum.	% Pres.
		DAIV	SO	Human	Agree.	Labels
MNIST	DX	35	35	81	<u>93</u>	92
	DLF	37	85	100	99	99
	SV	97	100	100	<u>96</u>	58
	FPT	<u>90</u>	100	<u>99</u>	<u>94</u>	54
	DJ	97	100	100	<u>96</u>	93
SVHN	DX	51	<u>99</u>	77	<u>68</u>	79
	DLF	100	100	96	<u>73</u>	50
	SV	100	100	<u>90</u>	<u>74</u>	9
	FPT	<u>99</u>	100	81	<u>75</u>	46
	DJ	0	1	61	76	9
ImageNet-1K	DX	100	100	<u>90</u>	100	<u>94</u>
	DLF	100	100	100	100	<u>95</u>
	SV	100	100	60	50	<u>83</u>
	FPT	100	100	100	100	100

EXTRA SLIDES

		TEST GENERATOR				
		DX	DLF	SV	FPT	DJ
DATASET	MNIST	 (a)	 (d)	 (g)	 (j)	 (m)
	SVHN	 (b)	 (e)	 (h)	 (k)	 (n)
	INet-1K	 (c)	 (f)	 (i)	 (l)	