

# CNN Shape vs Texture

Присяжнюк Артем

# Shape hypothesis: Интуиция

- Слой сопоставляет некоторой категории определенные формы
- Первые слои находят простые формы
- Средние слои находят сложные формы, состоящие из простых
- Последние слои находят формы объектов

# Texture hypothesis: Интуиция

Текстура влияет на решение слоя больше, чем форма



(a) Texture image

81.4%	<b>Indian elephant</b>
10.3%	indri
8.2%	black swan



(b) Content image

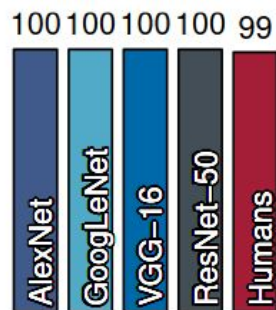
71.1%	<b>tabby cat</b>
17.3%	grey fox
3.3%	Siamese cat



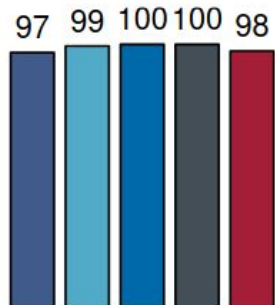
(c) Texture-shape cue conflict

63.9%	<b>Indian elephant</b>
26.4%	indri
9.6%	black swan

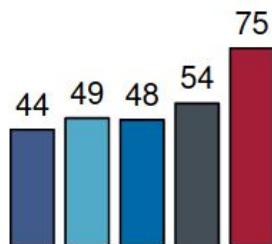
# Эксперимент: Датасет



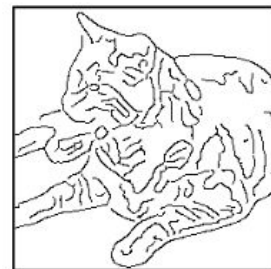
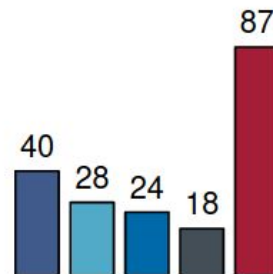
original



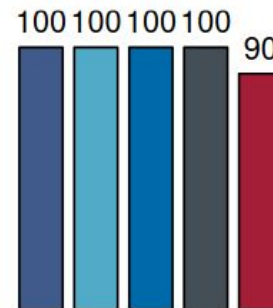
greyscale



silhouette

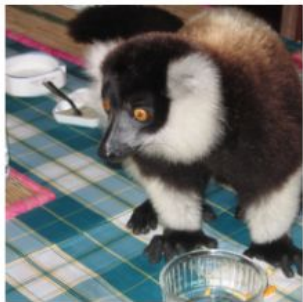


edges

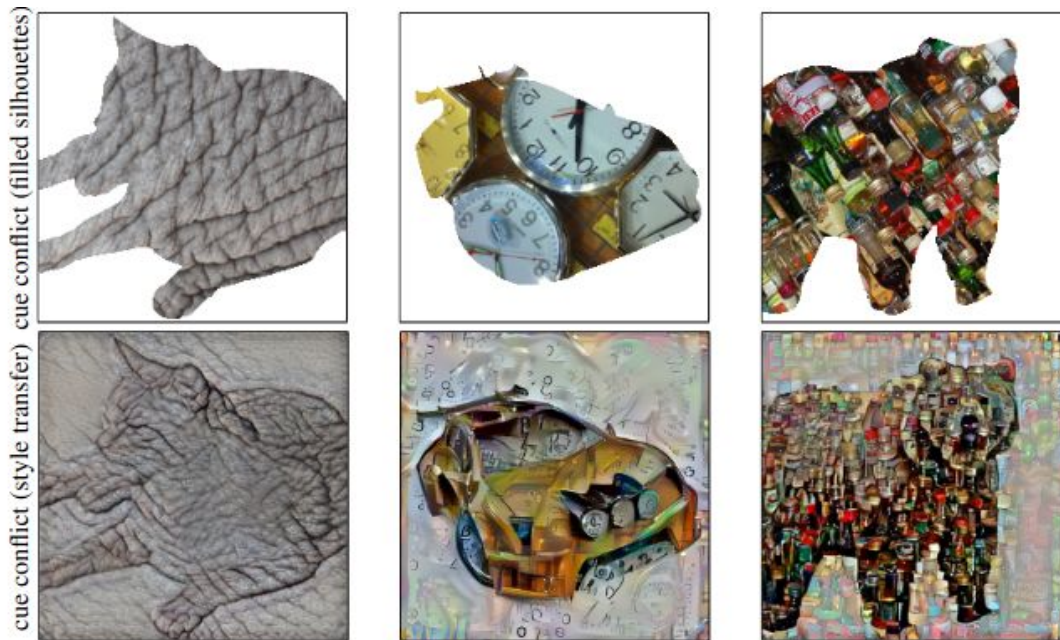


texture

# Эксперимент: Style Transfer



# Эксперимент: Stylized-ImageNet





# Эксперимент: Результаты

Красные - люди

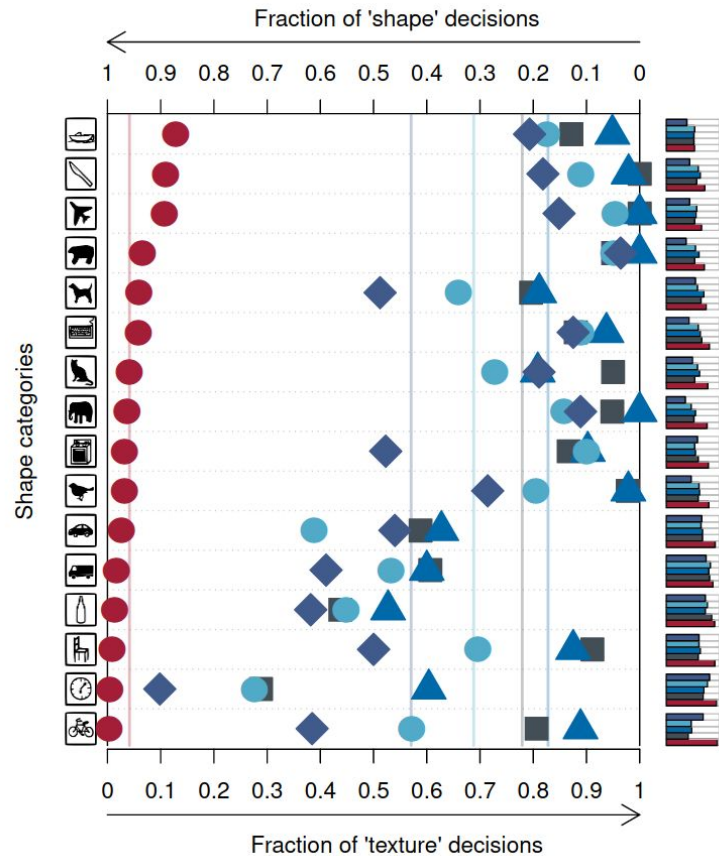
Фиолетовые - AlexNet

Синие - VGG-16

Бирюзовые - GoogLeNet

Серые - ResNet-50

architecture	IN→IN	IN→SIN	SIN→SIN	SIN→IN
ResNet-50	92.9	16.4	79.0	82.6
BagNet-33 (mod. ResNet-50)	86.4	4.2	48.9	53.0
BagNet-17 (mod. ResNet-50)	80.3	2.5	29.3	32.6
BagNet-9 (mod. ResNet-50)	70.0	1.4	10.0	10.9



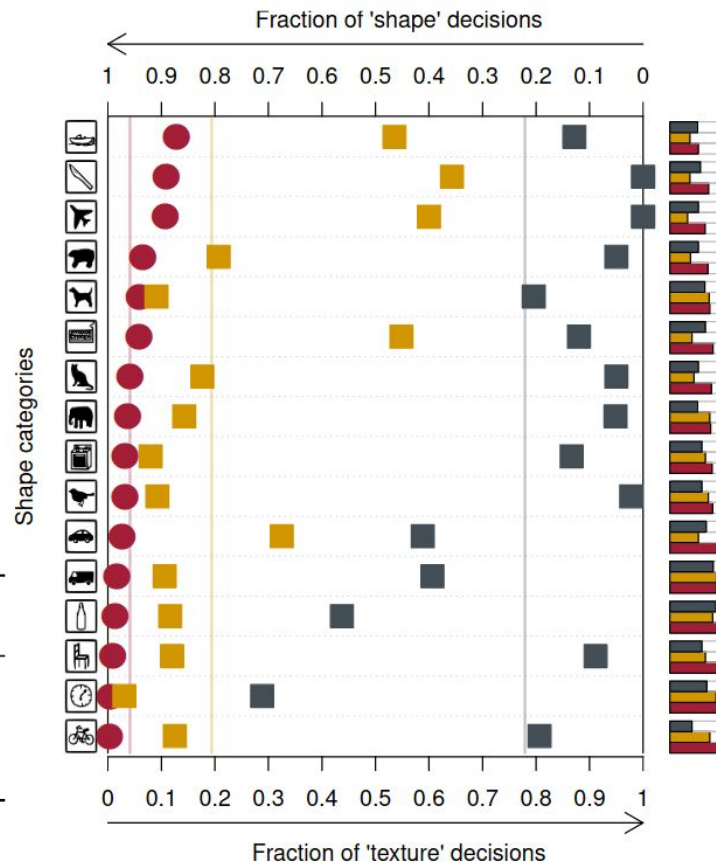
# Изменения: Результаты

Красные - люди

Золотые - Stylized-ImageNet

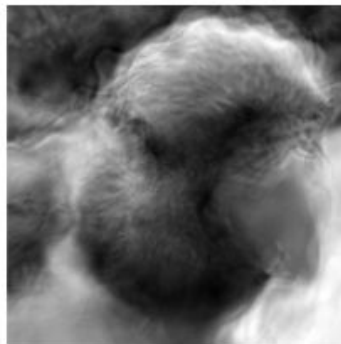
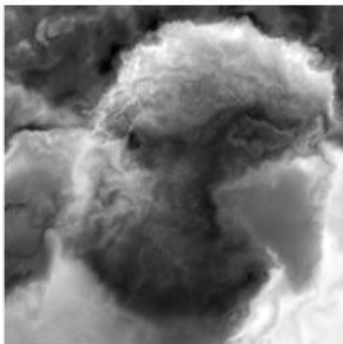
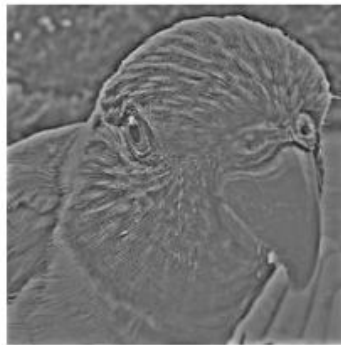
Серые - ImageNet

name	training	fine-tuning	top-1 IN accuracy (%)	top-5 IN accuracy (%)	Pascal VOC mAP50 (%)
vanilla ResNet	IN	-	76.13	92.86	70.7
	SIN	-	60.18	82.62	70.6
	SIN+IN	-	74.59	92.14	74.0
Shape-ResNet	SIN+IN	IN	<b>76.72</b>	<b>93.28</b>	<b>75.1</b>

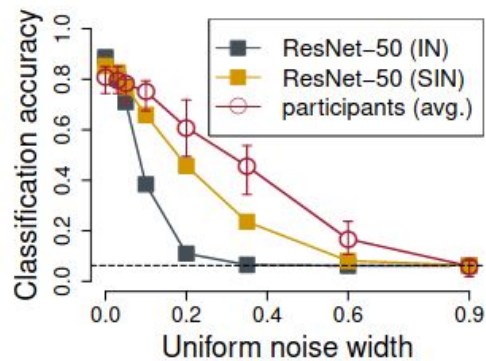




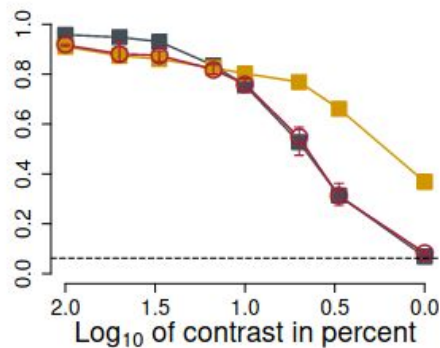
# Устойчивость к искажениям



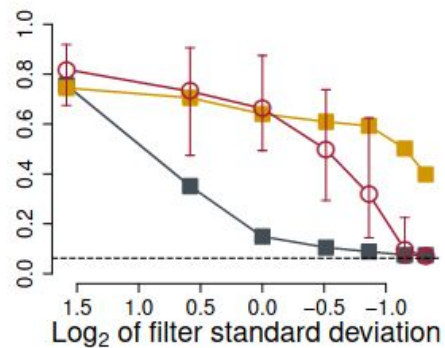
# Устойчивость к искажениям



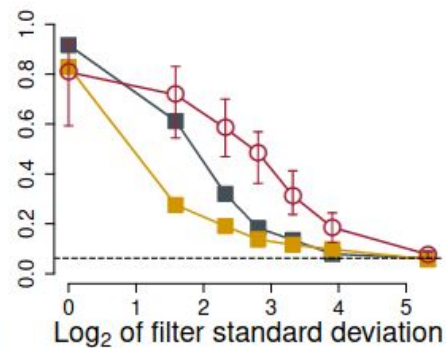
(a) Uniform noise



(b) Contrast



(c) High-pass



(d) Low-pass

# Итоги

- Обучение на ImageNet → texture bias
- Решение - Stylized-ImageNet (SIN)
- Обучение на SIN → устойчивость к искажениям
- Лучший вариант - Shape-ResNet
  - Обучение на SIN + IN
  - Fine-tuning на IN

# Сравнение моделей

CONVOLUTIONAL NEURAL NETWORKS			
MODEL	TOP-1(%)	MCE(%)	#PARAMS(M)
RESNET-50	76.02	65.54	26
ALEXNET	56.44	83.18	61
GOOGLENET	71.70	68.82	7
VGG-16	69.63	75.10	138

VISION TRANSFORMER ARCHITECTURES			
MODEL	TOP-1(%)	MCE(%)	#PARAMS(M)
ViT_BASE	75.73	58.55	86
ViT_LARGE	79.16	49.02	304
DeiT_BASE	81.84	42.30	86
DeiT_BASE-DIST.	83.16	41.19	87
DeiT_SMALL	79.68	47.79	22
DeiT_SMALL-DIST.	81.05	46.25	22
DeiT_TINY	71.92	<b>60.08</b>	<b>5</b>
DeiT_TINY-DIST.	74.38	57.45	6
CAiT_s24	83.28	40.59	47
CAiT_xxs24	78.38	49.28	11
SWIN-T_TINY	80.85	50.70	28
SWIN-T_SMALL	82.96	45.51	50
SWIN-T_BASE	84.90	38.52	88
SWIN-T_LARGE	85.92	<b>34.63</b>	<b>197</b>

CONVOLUTIONAL NEURAL NETWORKS		
MODEL	SHAPE BIAS (%)	# PARAMS (M)
RESNET-50	26.17	26
ALEXNET	29.80	61
GOOGLENET	28.52	7
VGG-16	16.12	138

VISION TRANSFORMER ARCHITECTURES		
MODEL	SHAPE BIAS (%)	# PARAMS (M)
ViT_BASE	49.10	86
ViT_LARGE	<b>55.35</b>	<b>304</b>
DeiT_BASE	42.32	86
DeiT_BASE-DIST.	39.62	87
DeiT_SMALL	38.26	22
DeiT_SMALL-DIST.	36.65	22
DeiT_TINY	<b>29.37</b>	<b>5</b>
DeiT_TINY-DIST.	31.06	6
CAiT_s24	38.65	47
CAiT_xxs24	34.24	11
SWIN-T_TINY	25.21	28
SWIN-T_SMALL	27.43	50
SWIN-T_BASE	36.39	88
SWIN-T_LARGE	40.20	197