

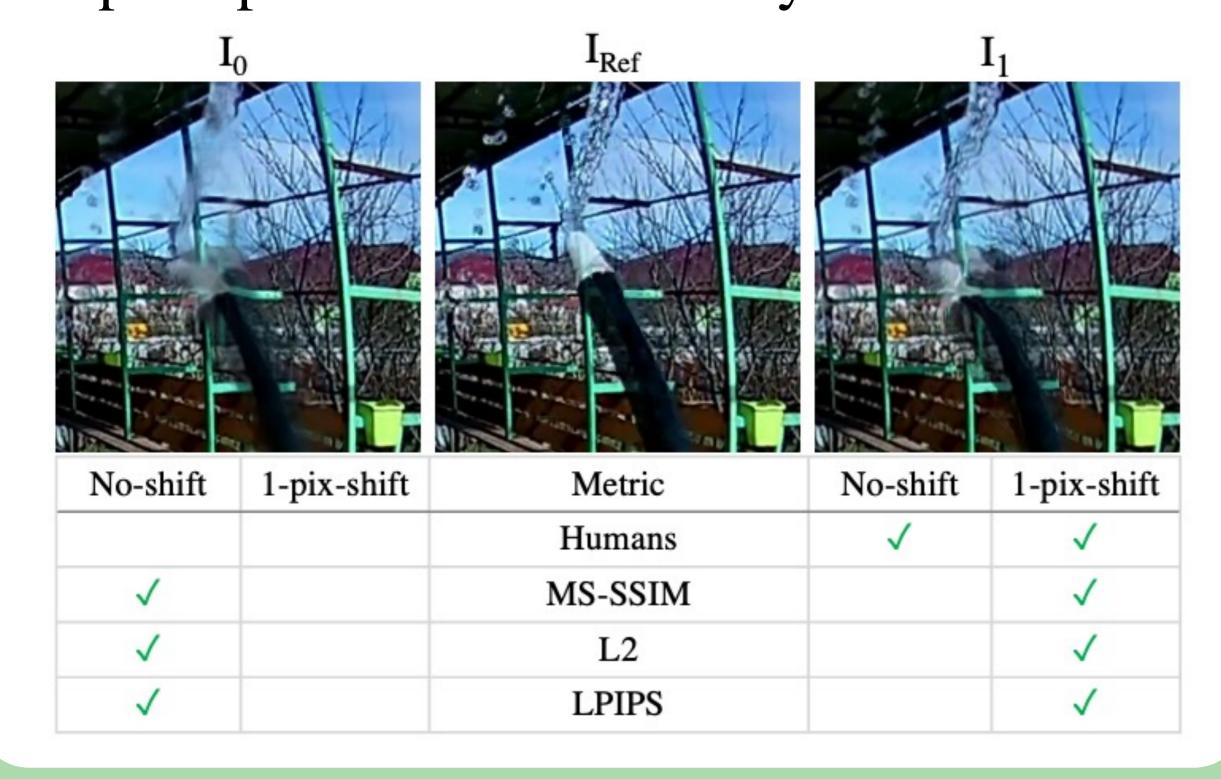
Shift-tolerant Perceptual Similarity Metric

Abhijay Ghildyal and Feng Liu



Overview

Existing perceptual similarity metrics assume an image and its reference are well aligned. As a result, these metrics are often sensitive to a small alignment error that is imperceptible to the human eyes.



Motivation

As commonly expected, shifting one image by a few pixels will not alter human similarity judgment on a pair of images. We conducted a user study to verify this common belief.

Pixel	Number	of user res	Avg. of std. in	
shift		Said No (Shifted)	Yes%	user responses per sample
0	140	10	93.3%	0.09 ± 0.17
1	121	29	80.7%	0.19 ± 0.23
2	84	66	56.0%	0.34 ± 0.21
3	52	98	34.7%	0.24 ± 0.23
4	52	98	34.7%	0.30 ± 0.24
5	40	110	26.7%	0.23 ± 0.24
6	35	115	23.3%	0.21 ± 0.24
7	31	119	20.7%	0.12 ± 0.20
8	27	123	18.0%	0.18 ± 0.23
9	15	135	10.0%	0.13 ± 0.21

Method

We build upon the Learned Perceptual Image Patch Similarity metric (LPIPS), a widely used learned metric, and explore architectural design considerations to make it robust against the imperceptible misalignment. Specifically, we study a wide spectrum of neural network elements, such as anti-aliasing filtering, pooling, striding, padding, and skip connection, and discuss their roles in making a robust metric.

Ablation Studies

We trained all our metrics using the original BAPPS training set on their original size of 256×256. We purposely did not train on the shifted version of the dataset to focus on neural

network element designs.

- Combining BlurPool with reducing stride size makes the network more robust against imperceptible shifts (r_{rf}) and more consistent with human judgment (2AFC).
- Blur Before Activation works better when the stride size 1.
- BlurPool significantly improves the robustness of other backbones as well. A larger padding size improves 2AFC.

AA (Bluz Reflectio	/	F-Conv	Stride	2AFC	·	r_{rf}	
1	2		in conv-1		1pixel	2pixel	3pixel
			4	70.65	2.87	3.92	3.74
\checkmark			2	70.53	1.85	2.22	2.58
	\checkmark		2	70.67	1.46	1.82	2.25
			4	70.57	2.78	3.92	3.91
	√	\checkmark	2	70.52	1.77	2.15	2.48
	<u> </u>		2	70.54	$\frac{1.84}{}$	$\frac{1}{2.28}$	$\frac{1}{2.34}$
	1		1	70.42	0.66	1.13	1.83
	\checkmark	\checkmark	1	70.44	0.63	1.14	1.68
§			2	70.57	2.63	3.36	3.16
	√§		2	70.63	2.80	3.57	3.39
	√ §	✓	2	70.52	2.95	4.13	3.93
Anti-Alia	ng C4-	ride	BlurPool	2AFC		r_{rf}	
(BlurPoo			Location	ZAFU		2pixel	2pivol

Anti-Alias	Stride	BlurPool	2AFC		3	
(BlurPool)	in $Conv-1$	Location		1pixel	2pixel	3pixe
√	2	Original	70.67	1.46	1.82	2.25
\checkmark	2	FeatAfterBlur	70.55	1.73	1.84	2.49
✓	2	$\underline{ BlurBeforeAct}$	70.50	2.06	2.02	2.74
<u>-</u>	1	Original	70.42	0.66	1.13	1.83
\checkmark	1	FeatAfterBlur	70.52	0.69	1.11	1.60
\checkmark	1	BlurBeforeAct	70.48	0.57	1.06	1.50

Network	`	lurPool) tion-Pad	2AFC		r_{rf}	
	1	2		1pixel	2pixel	3pixel
VGG-16			70.03	3.01	3.76	3.44
	\checkmark		70.05	0.66	1.08	1.44
		\checkmark	70.07	0.66	1.12	1.82
ResNet-18			69.86	$\frac{-}{2.67}$	3.35	$\frac{-}{3.77}$
	\checkmark		69.95	0.82	1.51	2.19
		\checkmark	70.14	1.07	1.81	2.38
Squeeze			$\frac{-}{69.61}$	$\frac{-}{7.41}$	$\frac{-7.58}{7.58}$	10.35
	\checkmark		69.24	2.03	3.06	3.93
		\checkmark	69.44	2.10	2.48	3.42

Results

Our metrics are both more robust against imperceptible shifts and consistent with human visual similarity judgment than most of the similarity metrics on BAPPS dataset.

Our method outperforms all other methods on the CLIC dataset while being shift-robust.

Based on the responses from the user study our metric follows the sensitivity of human perception to pixel shifts more accurately.

Network	2AFC		$rac{r_{rf}}{2 ext{pixel}}$	
		Tpixei	Zpixei	<u>spixei</u>
L2	62.92	3.59	7.55	10.82
SSIM [30]	61.41	3.16	7.20	13.73
CW-SSIM [31]	61.48	3.91	6.88	9.47
MS-SSIM [32]	62.54	2.22	5.83	10.66
PIEAPP Sparse [25]	64.20	2.83	3.19	3.81
PIEAPP Dense [25]	64.15	2.97	1.37	3.33
PIM-1 [3]	67.45	0.79	1.70	2.52
PIM-5 [3]	67.38	1.01	1.88	2.96
GTI-CNN [21]	63.87	3.95	4.91	7.88
DISTS [6]	68.83	2.85	2.89	4.03
E-LPIPS [16]	68.22	5.84	5.86	5.77
LPIPS (Alex) [37]	68.59	2.81	3.41	3.84
LPIPS (Alex) §*†	70.54	2.58	3.59	3.53
LPIPS (Alex) ours*†	70.39	0.66	1.24	1.79
LPIPS (Alex) §*‡	70.65	2.87	3.92	3.74
LPIPS (Alex) ours*‡	70.48	0.57	1.06	1 50

(§) Retrained from scratch. (*) Trained on patches of size 256 using author's (†) / our (‡) setup.

	. (~)	No. of rank flips			
Network	Accuracy(%)	$1_{\rm pixel}$	2pixel	3pixel	
L2	58.16	833	2102	2214	
SSIM [30]	60.00	349	931	1109	
PIEAPP [25]	75.44	91	134	158	
E-LPIPS [16]	74.44	212	251	317	
DISTS [6]	75.63	28	36	50	
PIM-1 [3]	73.79	13	22	33	
LPIPS(Alex) [37]	73.68	90	108	121	
$LPIPS(Alex)^{\S*\dagger}$	76.53	59	51	62	
LPIPS(Alex) ours*†	76.97	17	14	21	

(§) Retrained from scratch. (*) Trained on image

Metric	JND mAP%
SSIM [30]	0.722
LPIPS (Alex) [37]	0.757
LPIPS (Alex) §*†	0.740
LPIPS (Alex) ours*†	0.771
LPIPS (VGG) [37]	0.770
LPIPS (VGG) §*†	0.769
LPIPS (VGG) ours*†	0.775
DISTS [6]	0.766
PIM-1 [3]	0.773

(§) Retrained from scratch. (*) Trained on image patches of size 64 using author's (†) setup.

