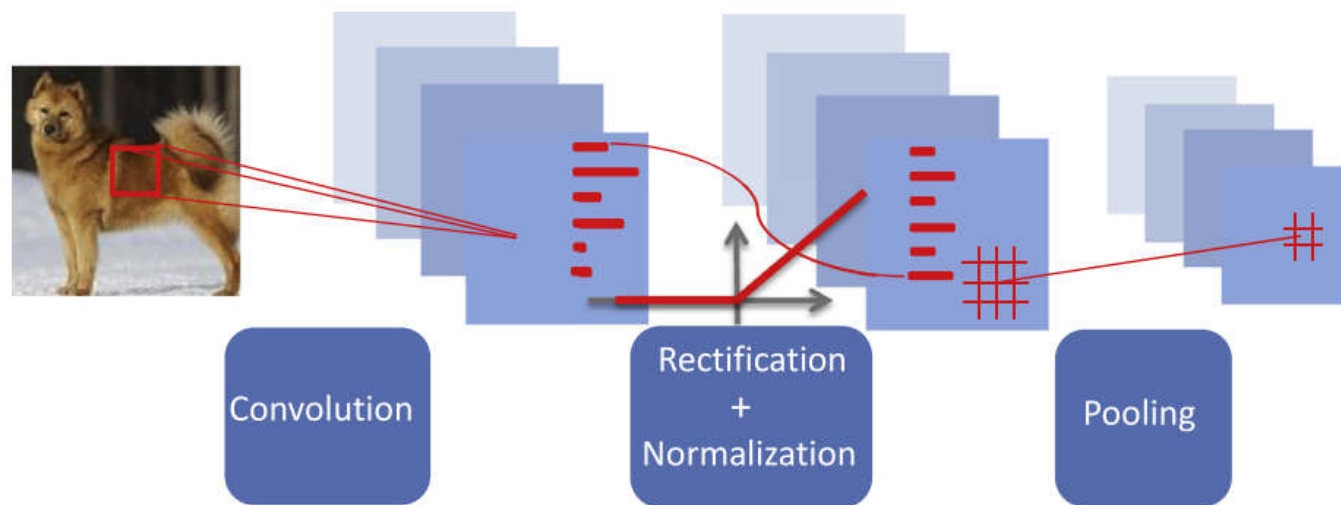# Detail-Preserving Pooling in Deep Networks

Zenglin Shi
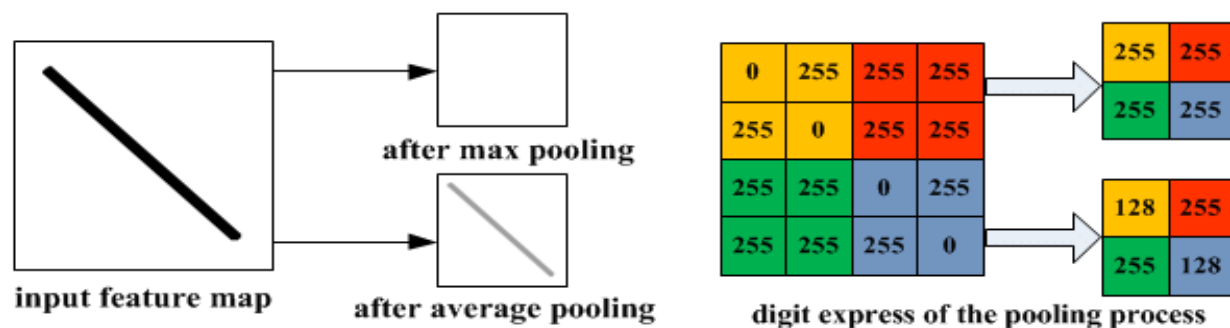
Dec. 10, 2018

# Pooling Mechanism
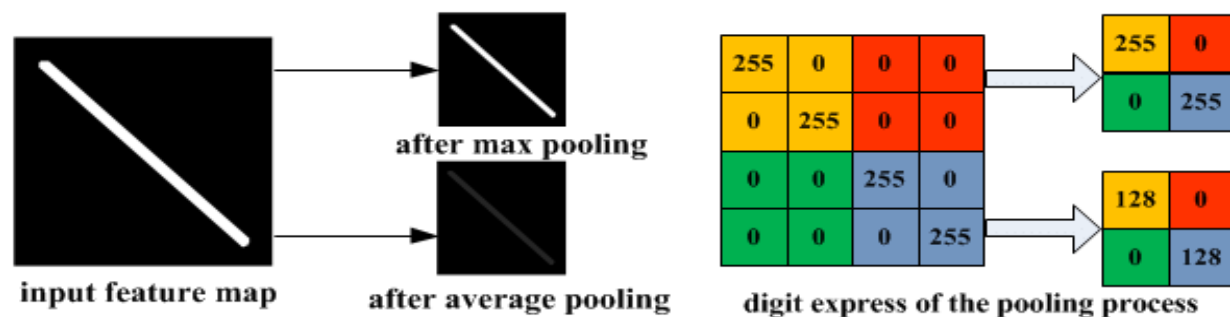


- Dimention Reduction
- Information Compression
- RF Enlarging
- Linear Transformation
- Translation Invariance

# Max & AVG. Pooling



after max pooling

input feature map

after average pooling

digit express of the pooling process

| 0 | 255 | 255 | 255 |
| 255 | 0 | 255 | 255 |
| 255 | 255 | 0 | 255 |
| 255 | 255 | 255 | 0 |

| 255 | 255 |
| 255 | 255 |

| 128 | 255 |
| 255 | 128 |

(a) Illustration of max pooling drawback

input feature map

after max pooling

after average pooling

digit express of the pooling process

| 255 | 0 | 0 | 0 |
| 0 | 255 | 0 | 0 |
| 0 | 0 | 255 | 0 |
| 0 | 0 | 0 | 255 |

| 255 | 0 |
| 0 | 255 |

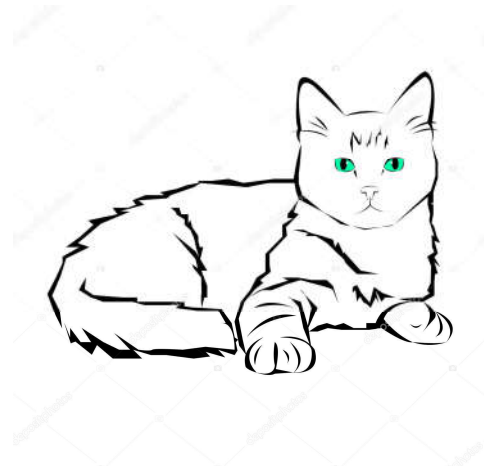| 128 | 0 |
| 0 | 128 |

(b) Illustration of average pooling drawback

# Weighted AVG. Pooling

# Detail-Preserving Pooling-Motivation
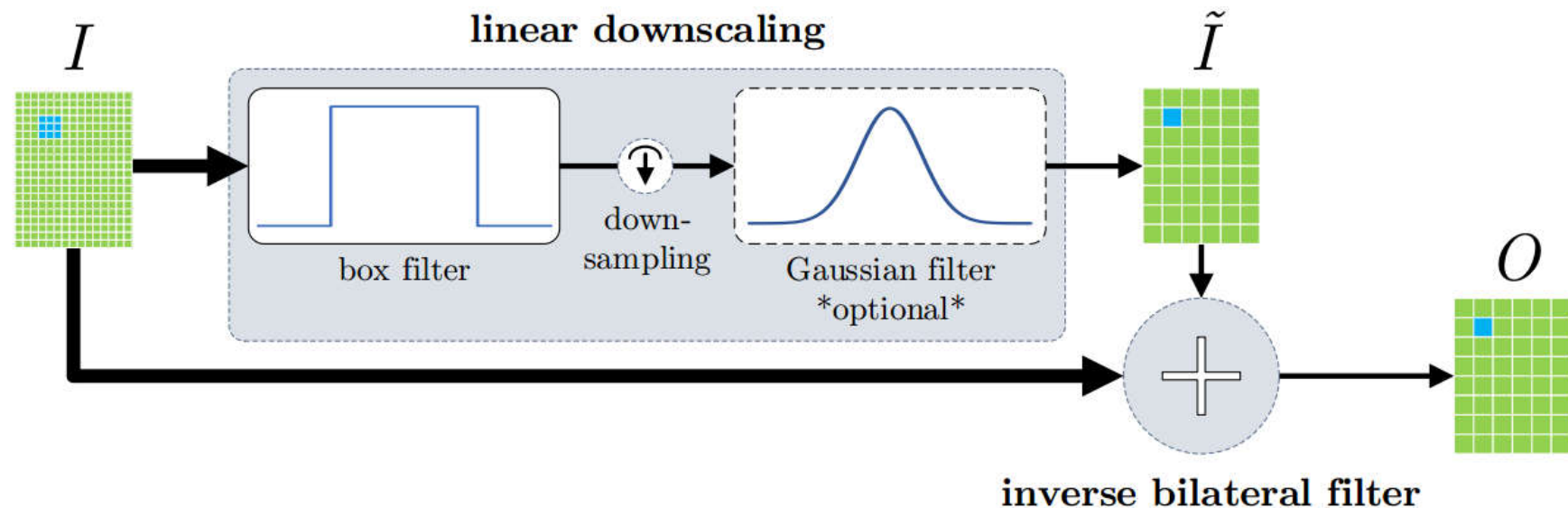
# Detail-Preserving Pooling-DPID



Figure 2. Diagram of detail-preserving downscaling (DPID) [31] and our detail-preserving pooling (DPP). DPP omits the Gaussian filter; Full-DPP replaces the box filter with a learned 2D filter.

# Detail-Preserving Pooling-DPID

- Given an input image I[·], detail-preserving image downscaling (DPID) calculates the downscaled output at pixel p as

$$O[p] = \frac{1}{k_p} \sum_{q \in \Omega_p} I[q] \cdot \left\| I[q] - \tilde{I}[p] \right\|^{\lambda}, \qquad (1)$$

- in which the linearly downscaled image I˜ is given by

$$\tilde{I} = I_D * \frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}^T \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}. \qquad (2)$$

- The weights are normalized with $k_p$

$$k_p = \sum_{q \in \Omega_p} \left\| I[q] - \tilde{I}[p] \right\|^{\lambda} \qquad (3)$$

# Detail-Preserving Pooling-DPID



(a) Original Image

(b) DPID ($\lambda = 1$)

(c) Box filter (average pooling)

(d) Spatial maximum \ extremum

# Detail-Preserving Pooling-DPP

- Detail-preserving pooling of an input activation map I at spatial output position p as

$$\mathcal{D}_{\alpha,\lambda}(I)[p] = \frac{1}{\sum_{q' \in \Omega_p} w_{\alpha,\lambda}[p,q']} \sum_{q \in \Omega_p} w_{\alpha,\lambda}[p,q] I[q]. \quad (4)$$

- Equation (4) computes a spatially weighted average of the input nodes in a neighborhood $I[q]_{q \in \Omega_p}$ for which we define weights $w_{\alpha,\lambda}[p, q]$ as

$$w_{\alpha,\lambda}[p,q] = \alpha + \rho_\lambda \left( I[q] - \tilde{I}[p] \right). \quad (5)$$

# Detail-Preserving Pooling-DPP

- For the symmetric variant of the reward function ρλ(·), employ the differentiable (generalized) Charbonnier penalty,

$$\rho_{\text{Sym}}(x) = \left(\sqrt{x^2 + \epsilon^2}\right)^{\lambda} \quad (6)$$

- The asymmetric variant of $\rho_{\lambda}(\cdot)$ only rewards positive arguments and is formulated as,

$$\rho_{\text{Asym}}(x) = \left(\sqrt{\max(0, x)^2 + \epsilon^2}\right)^{\lambda}, \quad (7)$$

# Detail-Preserving Pooling-DPP

- For the sake of simplicity in notation, we reformulate this such that the weights are normalized as

$$\tilde{w}_{\alpha,\lambda}[p,q] = \frac{w_{\alpha,\lambda}[p,q]}{\sum_{q' \in \Omega_p} w_{\alpha,\lambda}[p,q']} \quad (8)$$

- which allows us to write DPP as

$$\mathcal{D}_{\alpha,\lambda}(I)[p] = \sum_{q \in \Omega_p} \tilde{w}_{\alpha,\lambda}[p,q]I[q] \quad (9)$$

- In Eq. (4), I˜ is the result of a linear downscaling, achieving full flexibility with

$$\tilde{I}_F[p] = \sum_{q \in \tilde{\Omega}_p} F[q]I[q] \quad (10)$$

# Detail-Preserving Pooling-DPP



(a) Orig. feature map  (b) Max pooling  (c) Average pooling

(d) Extremum pooling  (e) Asymmetric DPP $(\lambda = 1)$  (f) Symmetric DPP $(\lambda = 2.7)$

Figure 4. Visualization of different pooling methods on an example feature map taken from the second layer of VGG-16. For both reward variants, the bias $\alpha$ is set to 0 to visually magnify the effect of the inverse bilateral weights. *Best viewed on screen.*

# Detail-Preserving Pooling-DPP

| | Method | VGG | NIN | ResNet |
|---|---|---|---|---|
| Deterministic methods | Strided conv. | $8.43 \pm 0.20$ | $10.97 \pm 0.10$ | $6.23^{(*)}$ |
| | Max | $7.43 \pm 0.20^{(*)}$ | $9.42 \pm 0.07$ | $6.52$ |
| | Average | $7.12 \pm 0.18$ | $8.75 \pm 0.15$ | $6.33$ |
| | NIN | – | $9.01 \pm 0.11^{(*)}$ | – |
| | Mixed (50/50) | $7.27 \pm 0.20$ | $8.68 \pm 0.23$ | $6.05$ |
| | Gated | $7.25 \pm 0.14$ | $8.67 \pm 0.22$ | $7.12$ |
| | $L_2$ | $7.15 \pm 0.18$ | $8.65 \pm 0.12$ | $7.29$ |
| | Lite-DPP$_{Asym}$ | $\underline{7.10} \pm 0.15$ | $8.62 \pm 0.10$ | $6.17$ |
| | Full-DPP$_{Asym}$ | $7.17 \pm 0.18$ | $8.73 \pm 0.05$ | $6.23$ |
| | Lite-DPP$_{Sym}$ | $7.19 \pm 0.10$ | $8.58 \pm 0.11$ | $6.05$ |
| | Full-DPP$_{Sym}$ | $\mathbf{7.02} \pm 0.18$ | $8.70 \pm 0.14$ | $5.97$ |
| Stoch. | Stochastic | $7.67 \pm 0.10$ | $8.92 \pm 0.09$ | $5.83$ |
| | S3pool | $7.21 \pm 0.14$ | $\underline{7.23} \pm 0.08$ | $\underline{5.55}$ |
| | Lite-S3DPP$_{Sym}$ | – | $\mathbf{7.13} \pm 0.09$ | $\mathbf{5.42}$ |

# Detail-Preserving Pooling in Deep Networks

Faraz Saeedan[1]  Nicolas Weber[1,2]  Michael Goesele[1,3]  Stefan Roth[1]

[1]TU Darmstadt  [2]NEC Laboratories Europe  [3]Oculus Research

CVPR 2018 — SALT LAKE CITY · JUNE 18-22

TECHNISCHE UNIVERSITÄT DARMSTADT

## Motivation

original — detail rich

spatial average — detail smoothed away

Similarly preserve detail, while pooling feature maps in CNNs

spatial max \ extremum — significant artefacts (noise) introduced

DPID [1] — detail plausibly preserved

## Detail-preserving image downscaling

- DPID [1]: downscale such that pixels that deviate more from a guidance $\tilde{I}$ have larger contribution

input $I$   guidance $\tilde{I}$   weights $\|I - \tilde{I}\|$

Similar to inverse bilateral filtering

$$O[p] = \frac{1}{k_p} \sum_{q \in \Omega_p} I[q] \cdot \|I[q] - \tilde{I}[p]\|^\lambda$$

output $O$

What is the role of $\lambda$?

Level of detail-preservation is adjusted by $\lambda$

## Detail-Preserving Pooling (DPP)

Using the intuition from DPID, we propose a pooling layer for CNNs that acts on every feature map independently

$$\mathcal{D}_{\alpha,\lambda}(I)[p] \propto \sum_{q \in \Omega_p} w_{\alpha,\lambda}[p,q] I[q]$$

- Weights are designed such that they magnify activations that stand out from their neighbors

- Larger bias $\alpha$ increases averaging tendency

reward function

$$w_{\alpha,\lambda}[p,q] = \alpha + \rho_\lambda \left( I[q] - \tilde{I}[p] \right)$$

bias

- Linear filtering before downsampling:
  $F$ can be a learnable 2D filter (**Full-DPP**) or a simple spatial average (**Lite-DPP**)

- Two choices for the reward function $\rho_\lambda(\cdot)$:

- **DPP_Sym**: In the limit as $\lambda \to \infty$, the symmetric reward models *extremum pooling*, i.e., it chooses the pixel with the largest deviation from the local average

- **DPP_Asym**: The asymmetric reward models classical *max pooling* in the limit

average pooling   extremum pooling

average pooling   max pooling

DPP is a generic pooling layer

$$\rho_{Sym}(x) = \left(\sqrt{x^2 + \epsilon^2}\right)^\lambda$$

$$\rho_{Asym}(x) = \left(\sqrt{\max(0, x)^2 + \epsilon^2}\right)^\lambda$$

## Properties

- DPP is **differentiable** w.r.t input and parameters
- For $\lambda = 0$, Lite-DPP equals *average pooling*
- For large values of $\lambda$, Lite-DPP_Sym (Asym) equals *extremum pooling* (max pooling)

## Parameters and complexity

- Lite-DPP has 2 parameters per feature map $(\lambda, \alpha)$
- Full-DPP has extra 10 parameters (3x3 filter) per feature map
  - Example: Full-DPP adds 0.172% and 0.098% parameters to ResNet-50 and 101
- Time for pooling is **independent of network depth**; quite minor for deeper networks
  - Worst case scenario in our experiments: 20% overhead for VGG-16

## Learnability and behavior adjustment

- Owing to differentiability:
  → Pooling behavior is adjusted through learning the parameters
- $\lambda$ and $\alpha$ are constrained to remain positive

DPP adjusts its behavior for every feature channel by learning all parameters

## Experiments: CIFAR-10

Three different networks, basic pooling methods and recent compound pooling layers compared

DPP outperforms other pooling methods for all networks

| Method | VGG | NiN | ResNet |
|---|---|---|---|
| Strided conv. | 8.43 | 10.97 | 6.23 |
| Max | 7.43 | 9.42 | 6.52 |
| Average | 7.12 | 8.75 | 6.33 |
| NiN | – | 9.41 | – |
| Mixed (NiN) | 7.27 | 8.68 | 6.02 |
| Gated | 7.58 | 8.67 | 7.12 |
| $L_2$ | 7.13 | 8.65 | 7.28 |
| Lite-DPP_Sym | 7.10 | 8.42 | 6.17 |
| Full-DPP_Sym | 7.17 | 8.71 | 6.22 |
| Lite-DPP_Asym | 7.18 | 8.40 | 5.99 |
| Full-DPP_Asym | **7.02** | **8.38** | **6.07** |

(*) indicates the original choice of pooling for the networks

## Experiments: ImageNet

Data augmentation and preparation according to [2]

DPP increases accuracy for all networks

| Network | Single crop | | Ten crop | |
|---|---|---|---|---|
| | Top-1 | Top-5 | Top-1 | Top-5 |
| VGG-16 (original) | 30.40 | 10.20 | – | – |
| VGG-16 + Average | 30.36 | 10.59 | – | – |
| VGG-16 + Lite-DPP_Sym | **30.05** | **10.47** | – | – |
| ResNet-50 | 24.25 | 7.36 | 22.33 | 6.20 |
| ResNet-50 + Lite-DPP_Sym | 23.42 | 6.83 | 21.82 | 5.92 |
| ResNet-101 | 22.35 | 6.23 | 20.82 | 5.38 |
| ResNet-101 | 22.16 | 6.16 | 20.89 | 5.21 |
| ResNet-101 + Lite-DPP_Sym | **21.79** | **5.91** | **20.13** | **5.20** |

## Joint regularization and pooling

- Similar to [3], S3DPP consists of DPP with a stride of 1, followed by a stochastic pixel selection
- Performs importance sampling
- S3DPP helps with both downscaling and regularization

| Method | VGG | NiN | ResNet |
|---|---|---|---|
| Strided conv. | 8.43 | 10.97 | 6.23 |
| Max | 7.43 | 9.42 | 6.52 |
| Average | 7.12 | 8.75 | 6.33 |
| NiN | – | 9.41 | – |
| Stochastic | 7.67 | 8.42 | 5.97 |
| S3pool | 7.25 | 7.23 | 5.19 |
| Lite-S3DPP_Sym | | **7.13** | **5.02** |

S3DPP combines stochastic regularization and pooling

## Learned parameters

- VGG-16 with 5 pooling layers trained on CIFAR-10
- Lite-DPP_Sym used for pooling
- Similar results for DPP_Asym

First layers act similar to max, later layers similar to average pooling

Similar to average pooling

Similar to max pooling

## Conclusion

- DPP is a learnable pooling layer that can perform similar to max / extremum or average pooling, or on a **nonlinear continuum** of intermediate functions
- Can be paired with regularization to further improve performance
- Constant and **affordable computational cost** and low parameter count
- **Consistently improves the performance** of networks on a wide range of architectures and datasets

Plug & play pooling layer

Code available online

## References

[1] N. Weber, M. Waechter, S. Amend, S. Guthe, and M. Goesele. Rapid, detail-preserving image downscaling. In *SIGGRAPH Asia 2016*.
[2] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR 2016*.
[3] S. Zhai, H. Wu, A. Kumar, Y. Cheng, Y. Lu, Z. Zhang, and R. Feris. S3pool: Pooling with stochastic spatial sampling. In *CVPR 2017*.