

6-1-2021

Why Dilated Convolutional Neural Networks: A Proof of Their Optimality

Jonatan Contreras

The University of Texas at El Paso, jmcontreras2@utep.edu

Martine Ceberio

The University of Texas at El Paso, mceberio@utep.edu

Vladik Kreinovich

The University of Texas at El Paso, vladik@utep.edu

Follow this and additional works at: https://scholarworks.utep.edu/cs_techrep



Part of the [Computer Sciences Commons](#), and the [Mathematics Commons](#)

Comments:

Technical Report: UTEP-CS-21-38c

Published in *Entropy*, 2021, Vol. 23, Paper 767.

Recommended Citation

Contreras, Jonatan; Ceberio, Martine; and Kreinovich, Vladik, "Why Dilated Convolutional Neural Networks: A Proof of Their Optimality" (2021). *Departmental Technical Reports (CS)*. 1571.

https://scholarworks.utep.edu/cs_techrep/1571

This Article is brought to you for free and open access by the Computer Science at ScholarWorks@UTEP. It has been accepted for inclusion in Departmental Technical Reports (CS) by an authorized administrator of ScholarWorks@UTEP. For more information, please contact lweber@utep.edu.

Why Dilated Convolutional Neural Networks: A Proof of Their Optimality

Jonatan Contreras, Martine Ceberio, and Vladik Kreinovich

University of Texas at El Paso, El Paso TX 79968, USA; jmcontreras2@utep.edu, mceberio@utep.edu, vladik@utep.edu

* Correspondence: vladik@utep.edu (V.K.)

Abstract: One of the most effective image processing techniques is the use of convolutional neural networks that use convolutional layers. In each such layer, the value of the layer’s output signal at each point is a combination of the layer’s input signals corresponding to several neighboring points. To improve the accuracy, researchers have developed a version of this technique, in which only data from *some* of the neighboring points is processed. It turns out that the most efficient case – called *dilated convolution* – is when we select the neighboring points whose differences in both coordinates are divisible by some constant ℓ . In this paper, we explain this empirical efficiency by proving that for all reasonable optimality criteria, dilated convolution is indeed better than possible alternatives.

Keywords: convolutional neural networks; dilated neural networks; optimality

1. Introduction

Convolutional layers: input and output. At present, one of the most efficient techniques in image processing and in other areas is a *convolutional neural network*; see, e.g., [1]. Convolutional neural networks include special types of layers that perform linear transformations.

Each such layer is characterized by integer-values parameters $\underline{X} \leq \bar{X}$, $\underline{Y} \leq \bar{Y}$, $d_{in} \geq 1$, and $d_{out} \geq 1$; then:

- the input to this layer consists of the values $F_{d'}(x', y')$, where d' , x' , and y' are integers for which $\underline{X} \leq x' \leq \bar{X}$, $\underline{Y} \leq y' \leq \bar{Y}$, and $1 \leq d' \leq d_{in}$; and
- the output of this layer consists of the values $G_d(x, y)$, where d , x , and y are integers for which $\underline{X} \leq x \leq \bar{X}$, $\underline{Y} \leq y \leq \bar{Y}$, and $1 \leq d \leq d_{out}$.

Convolutional layer: transformation. A general linear transformation takes the form

$$G_d(x, y) = \sum_{d'=1}^{d_{in}} \left(\sum_{x'=\underline{X}}^{\bar{X}} \sum_{y'=\underline{Y}}^{\bar{Y}} K_d(x, x', y, y', d') \cdot F_{d'}(x', y') \right), \quad (1)$$

for some coefficients $K_d(x, x', y, y', d')$.

Transformations performed by a convolutional layer are a specific case of such generic linear transformations, where the following two restrictions are imposed:

- first, each values $G_d(x, y)$ depends only on the values $F_{d'}(x', y')$ for which both differences $|x - x'|$ and $|y - y'|$ do not exceed some fixed integer L , and
- the coefficients $K_d(x, x', y, y', d')$ depend only on the differences $x - x'$ and $y - y'$:

$$K_d(x, x', y, y', d') = k_d(x - x', y - y', d') \quad (2)$$

for some coefficients $k_d(i, j, d')$ defined for all pairs (i, j) for which $|i|, |j| \leq L$.

Citation: Contreras, J.; Ceberio, M.; Kreinovich, V. Why Dilated Convolutional Neural Networks: A Proof of Their Optimality. *Entropy* **2021**, *1*, 0. <https://doi.org/>

Received:
Accepted:
Published:

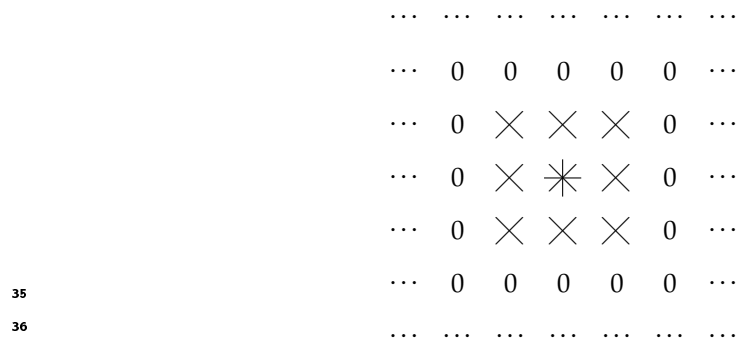
Publisher’s Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Copyright: © 2021 by the authors. Submitted to *Entropy* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

28 The values $k_d(i, j, d')$ are known as a *filter*.
 The resulting linear transformation takes the form

$$G_d(x, y) = \sum_{d'=1}^{d_{in}} \left(\sum_{-L \leq i, j \leq L} k_d(i, j, d') \cdot F_{d'}(x - i, y - j) \right). \quad (3)$$

29 Thus, the output $G_d(x, y)$ of a convolutional layer corresponding to the point (x, y)
 30 is determined by the values $F_{d'}(x - i, y - j)$ of the input to this layer at points $(x - i, y - j)$
 31 corresponding to $|i| \leq L$ and $|j| \leq L$. This is illustrated by Fig. 1, where, for $L = 1$ and
 32 for a point (x, y) marked by an asterisk, we show all the points $(x', y') = (x_0 - i, y_0 - j)$
 33 that determine the values $G_d(x, y)$. For convenience, points (x', y') that do not affect the
 34 values $G_d(x, y)$, are marked by zeros.



35
36
37
38

Figure 1: Convolution filter: case of $L = 1$

39 For $L = 2$, a similar picture has the following form:

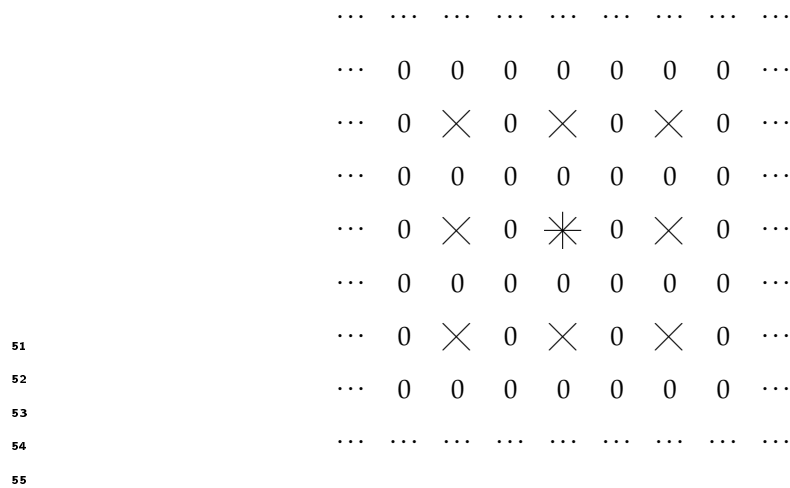


40
41
42
43

Figure 2: Convolution filter: case of $L = 2$

44 **Sparse filters and dilated convolution.** Originally, convolutional neural networks used
 45 filters in which all the values $k_d(i, j, d')$ for $|i|, |j| \leq L$ can be non-zero. It turned out,
 46 however, that we can achieve a better accuracy if we consider *sparse* filters, i.e., filters
 47 in which, for some pairs (i, j) with $|i|, |j| \leq L$, all the values $k_d(i, j, d')$ are fixed at 0; see,
 48 e.g., [3,5,6].

49 In Fig. 3, we show an example of such a situation, when $L = 2$ and only values
 50 $k_d(i, j, d')$ for which both i and j are even are allowed to be non-zero.



51
52
53
54
55
56
57 Figure 3. Case when $L = 2$ and only values $k_d(i, j, d')$ with even i and j can be non-zero

In general, it turned out that such a restriction works best if we only allow $k_d(i, j, d') \neq 0$ for pairs (i, j) which are divisible by some integer ℓ , i.e., if we take

$$G_d(x, y) = \sum_{d'=1}^{d=d_{in}} \left(\sum_{-L \leq i, j \leq L: i/\ell \in \mathbb{Z}, j/\ell \in \mathbb{Z}} k_d(i, j, d') \cdot F_{d'}(x - i, y - j) \right). \quad (4)$$

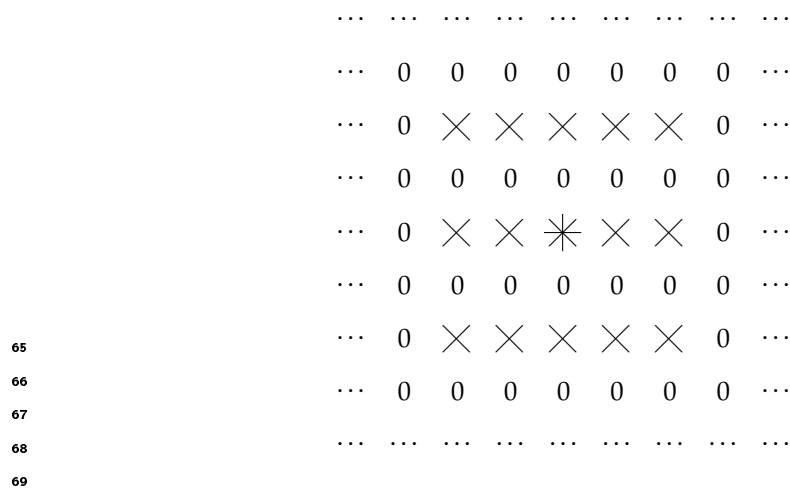
In this case, the layer's output signal $G_d(x, y)$ can be written in the following equivalent form:

$$G_d(x, y) = \sum_{d'=1}^{d_{in}} \left(\sum_{-\tilde{L} \leq \tilde{i}, \tilde{j} \leq \tilde{L}} \tilde{k}_d(\tilde{i}, \tilde{j}, d') \cdot F_{d'}(x - \ell \cdot \tilde{i}, y - \ell \cdot \tilde{j}) \right), \quad (5)$$

58 where we denoted $\tilde{L} \stackrel{\text{def}}{=} L/\ell$, $\tilde{i} \stackrel{\text{def}}{=} i/\ell$, $\tilde{j} \stackrel{\text{def}}{=} j/\ell$, and $\tilde{k}_d(\tilde{i}, \tilde{j}, d') \stackrel{\text{def}}{=} k(\ell \cdot \tilde{i}, \ell \cdot \tilde{j}, d')$.

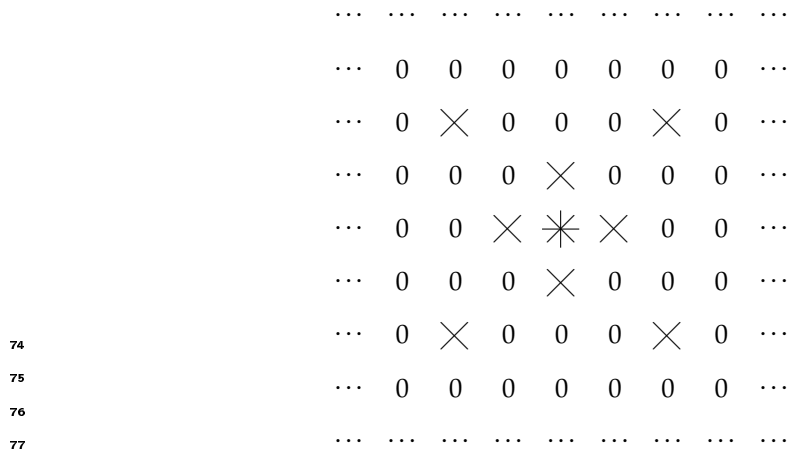
59 The resulting networks are known as *dilated* convolutional neural networks, since
60 skipping some points (i, j) in the description of the filter is kind of equivalent to extend-
61 ing (dilating) the distance between the remaining points; see, e.g., [3,5,6].

62 **Empirical fact that needs explanation.** In principle, we could select other points (i, j)
63 at which the filter can be non-zero. For example, we could select points for which j is
64 even, but i can be any integer:



65
66
67
68
69
70
71 Figure 4. Case when $L = 2$ and only values $k_d(i, j, d')$ with even j can be non-zero

72 Alternatively, for $L = 2$, as points (i, j) at which $k_d(i, j, d')$ can be non-zero, we
 73 could select the points $(0, 0)$, $(0, \pm 1)$, $(\pm 1, 0)$, and $(\pm 2, \pm 2)$, see Fig. 5.



80 Figure 5. A possible selection of points (i, j) for which $k_d(i, j, d')$ can be no-zero

81 However, empirical evidence shows that the selection corresponding to dilated
 82 convolution – when we select points for which i and j are both divisible by some integer
 83 ℓ – works the best [3,5,6].

84 To the best of our knowledge, there is *no theoretical explanation* for this empirical
 85 result – that dilated convolution leads to better results than selecting other sets of non-
 86 zero-valued points (i, j) . The main *objective* of this paper is to *provide* such an *explanation*.
 87

88 *Comment.* Let us emphasize that the only objective of this paper is to *explain* this em-
 89 pirical fact, we are not yet at a stage where we can propose a new method or even any
 90 improvements to the known methods.

91 **2. Analysis of the Problem**

Let us reformulate this situation in geometric terms: case of traditional convolution.

In the original convolution formula (1), to find the values $G_d(x, y)$ the layer’s output signal at a point (x, y) , we need to know the values $F_{d'}(x', y')$ the layer’s input signal at all the points (x', y') of the type $(x - i, y - j)$ for $|i|, |j| \leq L$. We can reformulate it by saying that we need to know the values $F_{d'}(x', y')$ at all the points (x', y') at which the ℓ_∞ distance

$$d_\infty((x, y), (x', y')) \stackrel{\text{def}}{=} \max(|x - x'|, |y - y'|), \tag{6}$$

does not exceed L :

$$G_d(x, y) = \sum_{d'=1}^{d_{\text{in}}} \left(\sum_{(x', y') \in D: d_\infty((x, y), (x', y')) \leq L} k_d(x - x', y - y', d') \cdot F_{d'}(x', y') \right), \tag{7}$$

where we denoted

$$D \stackrel{\text{def}}{=} (\mathbb{Z} \cap [\underline{X}, \bar{X}]) \times (\mathbb{Z} \cap [\underline{Y}, \bar{Y}]). \tag{8}$$

That we use, in this formula, the bounded subset D of the “grid” $\mathbb{Z} \times \mathbb{Z}$ and not the whole set $\tilde{S} \stackrel{\text{def}}{=} \mathbb{Z} \times \mathbb{Z}$ only matters at the border of the domain D . So, to simplify our

formulas, we can follow the usual tradition (see, e.g., [5]) and simply use the whole set $\tilde{S} = \mathbb{Z} \times \mathbb{Z}$ instead of the bounded set D :

$$G_d(x, y) = \sum_{d'=1}^{d_{\text{in}}} \left(\sum_{(x', y') \in \tilde{S}: d_{\infty}((x, y), (x', y')) \leq L} k_d(x - x', y - y', d') \cdot F_{d'}(x', y') \right). \quad (9)$$

92

93 *Comment.* Note that the set \tilde{S} is potentially *infinite*. What makes the set of all the points
 94 (x', y') – that affects the values $G_d(x, y)$ – *finite* is the restriction $d_{\infty}((x, y), (x', y')) \leq L$,
 95 whose meaning is that such points (x', y') should belong to the corresponding neighbor-
 96 hood of the point (x, y) .

Case of dilated convolution. The dilated convolution can be described in a similar way. Namely, we can describe the formula (4) as

$$G_d(x, y) = \sum_{d'=1}^{d_{\text{in}}} \left(\sum_{(x', y') \in S_{\ell}(x, y): d_{\infty}((x, y), (x', y')) \leq L} k_d(x - x', y - y', d') \cdot F_{d'}(x', y') \right), \quad (10)$$

where $S_{\ell}(x, y)$ denotes the set of all the points (x', y') for which both differences $x - x'$ and $y - y'$ are divisible by ℓ :

$$S_{\ell}(x, y) \stackrel{\text{def}}{=} \{(x', y') : x' \equiv x \pmod{\ell}, y' \equiv y \pmod{\ell}\}. \quad (11)$$

97 Note that in this representation of dilated convolution, while we have *several* dif-
 98 ferent sets $S_{\ell}(x, y)$ for different points (x, y) , there is only *one* filter $k_d(x - x', y - y', d')$,
 99 namely the same filter that was used in the original representation (4). So, in this new
 100 representation, we have exactly as many parameters as before.

101 The main difference between this formula and the formula (9) is that, in contrast to
 102 the usual convolution (9), where the same set $\tilde{S} = \mathbb{Z} \times \mathbb{Z}$ could be used for all the points
 103 (x, y) , here, in general, we may need different sets $S_{\ell}(x, y)$ for different points (x, y) .

104 For example, if $\ell = 2$, then we need four such sets:

- for points (x, y) for which both x and y are even, the formula (10) holds for

$$S_2(0, 0) = S_2(0, 2) = \dots = S_{0,0}^{(\ell=2)} \stackrel{\text{def}}{=} \{(x, y) \in \mathbb{Z} \times \mathbb{Z} : x \text{ and } y \text{ are even}\}; \quad (12)$$

- 105 • for points (x, y) for which x is even but y is odd, the formula (10) holds for

$$S_2(0, 1) = S_2(0, 3) = \dots = S_{0,1}^{(\ell=2)} \stackrel{\text{def}}{=} \{(x, y) \in \mathbb{Z} \times \mathbb{Z} : x \text{ is even and } y \text{ is odd}\}; \quad (13)$$

- 106 • for points (x, y) for which x is odd but y is even, the formula (10) holds for

$$S_2(1, 0) = S_2(1, 2) = \dots = S_{1,0}^{(\ell=2)} \stackrel{\text{def}}{=} \{(x, y) \in \mathbb{Z} \times \mathbb{Z} : x \text{ is odd and } y \text{ is even}\}; \quad (14)$$

- finally, for points (x, y) for which x and y are both odd, the formula (10) holds for

$$S_2(1, 1) = S_2(1, 3) = \dots = S_{1,1}^{(\ell=2)} \stackrel{\text{def}}{=} \{(x, y) \in \mathbb{Z} \times \mathbb{Z} : x \text{ and } y \text{ are odd}\}. \quad (15)$$

In this case, instead of the single set $S_1(x, y) = \tilde{S}$ (s in the case of the traditional convolution), we have a set of such sets

$$\mathcal{F} = \{S_{0,0}^{(\ell=2)}, S_{0,1}^{(\ell=2)}, S_{1,0}^{(\ell=2)}, S_{1,1}^{(\ell=2)}\}. \quad (16)$$

To avoid confusion, we will call subsets of the original “grid” $\mathbb{Z} \times \mathbb{Z}$ sets, while the set of such sets will be called a *family*. In these terms, the formula (10) can be described as follows:

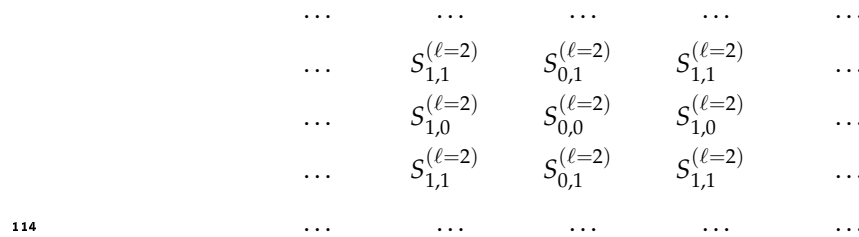
$$G_d(x, y) = \sum_{d'=1}^{d_{\text{in}}} \left(\sum_{(x', y') \in S_{\mathcal{F}}(x, y): d_{\infty}((x, y), (x', y')) \leq L} k_d(x - x', y - y', d') \cdot F_{d'}(x', y') \right), \quad (17)$$

where $S_{\mathcal{F}}(x, y)$ denotes the set $S \in \mathcal{F}$ from the family \mathcal{F} that contains the point (x, y) :

$$(x, y) \in S_{\mathcal{F}}(x, y) \text{ and } S_{\mathcal{F}}(x, y) \in \mathcal{F}. \quad (18)$$

107 In this representation, all four sets S from the family \mathcal{F} are *infinite* – just like the set
 108 \tilde{S} corresponding to the traditional convolution is infinite. Similarly to the traditional
 109 convolution, what makes the set of all the points (x', y') – that affects the values $G_r(x, y)$
 110 – *finite* is the restriction $d_{\infty}((x, y), (x', y')) \leq L$, whose meaning is that such points (x', y')
 111 should belong to the corresponding neighborhood of the point (x, y) .

112 Fig. 6 describes which of the four sets $S \in \mathcal{F}$ corresponds to each point (x, y) from
 113 the “grid” $\mathbb{Z} \times \mathbb{Z}$:



114

115

116

117

118

Figure 6. Sets $S_{\mathcal{F}}(x, y)$ corresponding to different points (x, y)
 – for filters presented in Figure 3

For $\ell = 3$, we can get a similar reformulation, with the family

$$\mathcal{F} = \left\{ S_{0,0}^{(\ell=3)}, S_{0,1}^{(\ell=3)}, S_{0,2}^{(\ell=3)}, S_{1,0}^{(\ell=3)}, S_{1,1}^{(\ell=3)}, S_{1,2}^{(\ell=3)}, S_{2,0}^{(\ell=3)}, S_{2,1}^{(\ell=3)}, S_{2,2}^{(\ell=3)} \right\}, \quad (19)$$

119 where $S_{i,j}^{(\ell=3)}$ is the set of all the pairs $(x, y) \in \mathbb{Z} \times \mathbb{Z}$ in which both differences $x - i$ and
 120 $y - j$ are divisible by 3.

121 In general, for an arbitrary point (x, y) , we should use the set $S_{\mathcal{F}} = S_{x \bmod \ell, y \bmod \ell}^{(\ell=2)}$.
 122

Other cases. Such a representation is possible not only for dilated convolution. For example, the above case when we allow arbitrary value i and require the value j to be even can be described in a similar way, with

$$\mathcal{F} = \{S_0, S_1\}, \quad (20)$$

123 where:

- for points (x, y) for which y is even, we take

$$S_{\mathcal{F}}(0, 0) = S_{\mathcal{F}}(1, 0) = \dots = S_0 \stackrel{\text{def}}{=} \{(x, y) \in \mathbb{Z} \times \mathbb{Z} : y \text{ is even}\}, \quad (21)$$

- and for points (x, y) for which y is odd, we take

$$S_{\mathcal{F}}(0, 1) = S_{\mathcal{F}}(1, 1) = \dots = S_1 \stackrel{\text{def}}{=} \{(x, y) \in \mathbb{Z} \times \mathbb{Z} : y \text{ is odd}\}. \quad (22)$$

124 In principle, we can also have families that have infinite number of sets; an example of
125 such a family will be given below.

126 We can also, in principle, consider the situations when we do not require that the
127 coefficients $k_d(x, x', y, y', d')$ depend only on the differences $x - x'$ and $y - y'$. Thus, we
128 arrive to the following general description.

129 **General case.** In the general case, we get the following situation:

- 130 • we have a family \mathcal{F} of subsets of the “grid” $\mathbb{Z} \times \mathbb{Z}$;
- 131 • the values $G_d(x, y)$ of the layer’s output signal at a point (x, y) is determined by the
132 formula

$$G_d(x, y) = \sum_{d'=1}^{d_{\text{in}}} \left(\sum_{(x', y') \in S_{\mathcal{F}}(x, y): d_{\infty}((x, y), (x', y')) \leq L} K_d(x, x', y, y', d') \cdot F_{d'}(x', y') \right), \quad (23)$$

133 for some values $K_d(x, x', y, y', d')$, where $S_{\mathcal{F}}(x, y)$ denotes the set $S \in \mathcal{F}$ from the
134 family \mathcal{F} that contains the point (x, y) .

135 For the formula (23) to uniquely determine the values $G_d(x, y)$, we need to make sure
136 that the set $S_{\mathcal{F}}(x, y)$ is uniquely determined by the point (x, y) , i.e., that for each point
137 (x, y) , the family \mathcal{F} contain one and only one set S that contains this point. In other
138 words:

- 139 • different sets from the family \mathcal{F} must be disjoint, and
- 140 • the union of all the sets $S \in \mathcal{F}$ must coincide with the whole “grid” $\mathbb{Z} \times \mathbb{Z}$.

141 In mathematical terms, the family \mathcal{F} must form a *partition* of the “grid” $\mathbb{Z} \times \mathbb{Z}$.

142 *Comment.* To avoid possible confusion, it is worth mentioning that while *different sets* S
143 from the family \mathcal{F} are disjoint, this does not preclude the possibility that *sets* $S_{\mathcal{F}}(x, y)$
144 and $S_{\mathcal{F}}(x', y')$ corresponding to *different points* $(x, y) \neq (x', y')$ can be identical. For
145 example, in the description of the traditional convolution, the family \mathcal{F} consists of only
146 one set $\mathcal{F} = \{\tilde{S}\}$. In this case, for all points (x, y) and (x', y') , we have $S_{\mathcal{F}}(x, y) =$
147 $S_{\mathcal{F}}(x', y') = \tilde{S}$.

148 In terms of sets corresponding to different points, disjointness means that *if* the
149 sets $S_{\mathcal{F}}(x, y)$ and $S_{\mathcal{F}}(x', y')$ are different, *then* these sets must be disjoint: $S_{\mathcal{F}}(x, y) \cap$
150 $S_{\mathcal{F}}(x', y') = \emptyset$.

151 **We do not a priori require shift-invariance.** Please note that we do not a priori require
152 that the sets $S_{\mathcal{F}}(x, y)$ and $S_{\mathcal{F}}(x_0, y_0)$ corresponding to two different points (x, y) and
153 (x_0, y_0) should be obtained from each other by shift – this property is known as *shift*
154 *invariance* and is satisfied both for the usual convolution and for the dilated convolution.

155 It should be emphasized, however, that we will show that this shift-invariance
156 property holds for the optimal arrangement.

157 **Let us avoid the degenerate case.** From the purely mathematical viewpoint, we can
158 have a partition of the “grid” $\mathbb{Z} \times \mathbb{Z}$ into one-point sets $\{(x, y)\}$. This is an example
159 when the family \mathcal{F} has infinitely many subsets.

160 In this case, no matter what value L we choose, the formula (23) implies that the
161 values $G_d(x, y)$ of the layer’s output signal at a point (x, y) is determined only by the
162 values $F_{d'}(x, y)$ of the layer’s input at this same point. This is equivalent to using a
163 convolution with $L = 0$; such a convolution is known as the 1-by-1 convolution.

164 While such convolution is often useful, in this case, for each point (x, y) , there is
165 only one point $(x', y') = (x, y)$, so it is not possible to select only *some* of the points
166 (x', y') – which is the whole idea of dilation. Since in this paper, we study dilation, we
167 will therefore avoid this 1-by-1 situation and additionally require that at least one set
168 from the family \mathcal{F} should contain more than one element.

169 **What we plan to do.** We will consider all possible families \mathcal{F} that form a partition of
170 the “grid” $\mathbb{Z} \times \mathbb{Z}$, and we will show that for all optimality criteria that satisfy some

171 reasonable conditions, the optimal family is either the family of sets corresponding to
 172 the dilated convolution – or a natural modification of this family.

173 Let us describe what we mean by an optimality criteria.

174 **What does “optimal” mean?** In our case, we select between different families of sets \mathcal{F} ,
 175 \mathcal{F}' , ... In general, we select between alternatives a , b , etc. Out of all possible alternatives,
 176 we want to select an *optimal* one. What does “optimal” mean?

177 In many cases, “optimal” is easy to describe:

- 178 • we have an objective function $f(a)$ that assigns a numerical value to each alternative
 179 a – e.g., the average approximation error of the numerical method a for solving a
 180 system of differential equations, and
- 181 • optimal means we select an alternative for which the value of this objective function
 182 is the smallest possible (or, for some objective functions, the largest possible).

183 However, this is not the only possible way to describe optimality.

184 For example, if we are minimizing the average approximation error, and there
 185 are several different numerical methods with the exact same smallest value of average
 186 approximation error, then we can use this non-uniqueness to select, e.g., the method with
 187 the shortest average computation time. In this case, we have, in effect, a more complex
 188 preference relation between alternatives than in the case when decision is made based
 189 solely on the value of the objective function. Specifically, in this case, an alternative b is
 190 better than the alternative a – we will denote it by $a < b$ – if:

- 191 • either we have $f(b) < f(a)$,
- 192 • or we have $f(a) = f(b)$ and $g(b) < g(a)$.

193 If this still leaves several alternatives which are equally good, then we can optimize
 194 something else and thus, have an even more complex optimality criterion.

195 In general, having an optimality criterion means that we are able to compare pairs
 196 of alternatives – at least some such pairs – and conclude that:

- 197 • for some of these pairs, we have $a < b$,
- 198 • for some of these pairs, we have $b < a$, and
- 199 • for some others pairs, we conclude that alternatives a and b are, from our viewpoint,
 200 of equal value; we will denote this by $a \sim b$.

201 Of course, these relations must satisfy some reasonable properties. For example, if b is
 202 better than a , and c is better than b , then c should be better than a ; in mathematical terms,
 203 the relation $<$ must be *transitive*.

204 What we *must* have is some alternative which is better than or equivalent to all
 205 others – otherwise, the optimization problem has no solutions. It also makes sense to
 206 require that there is only one such optimal alternative – indeed, as we have mentioned, if
 207 there are several equally good optimal alternatives, this means that the original optimal-
 208 ity criterion is not final, that we can use this non-uniqueness to optimize something else,
 209 i.e., in effect, to modify the original criterion into a final (or at least “more final”) one.

210 **Invariance.** There is an additional natural requirement for possible optimality criteria,
 211 which is related to the fact that the original “grid” $\mathbb{Z} \times \mathbb{Z}$ has lots of *symmetries*, i.e.,
 212 transformations that transform this “grid” into itself.

213 For example, if we change the starting point of the coordinate system to a new
 214 point (x_0, y_0) , then a point that originally had coordinates (x, y) now has coordinates
 215 $(x - x_0, y - y_0)$. It makes sense to require that the relative quality of two different families
 216 \mathcal{F} and \mathcal{F}' will not change if we simply change the starting point.

217 Similarly, we can change the direction of the x -axis, then a point (x, y) becomes
 218 $(-x, y)$. If we change the direction of the y -axis, we get a transformation $(x, y) \rightarrow (x, -y)$.
 219 Finally, we can rename the coordinates: what was x will become y and vice versa; this
 220 corresponds to the transformation $(x, y) \rightarrow (y, x)$. Such transformations should also not
 221 affect the relative quality of different families.

222 *Comment.* Please note that we are *not* requiring that the *family* \mathcal{F} of sets be shift-invariant,
 223 what we require is that the *optimality criterion* is shift-invariant.

224 Let us explain why, in our opinion, it makes sense to require that the optimality
 225 criterion is shift-invariance – as well as has other invariance properties. Indeed, let us
 226 consider any usual optimality criterion such as accuracy of classification, robustness to
 227 noise, etc. What each criterion means is, e.g., the overall classification accuracy over the
 228 set \mathcal{S} of all possible cat and not-a-cat images $I \in \mathcal{S}$. We want this method to correctly
 229 classify images into cats and not-cats, whether these images are centered or somewhat
 230 shifted. Thus, to adequately compare different methods, we should test these methods
 231 on a set \mathcal{S} of images that includes both original and shifted images.

232 Here:

- 233 • if we shift each image I from the set \mathcal{S} by the same shift (x_0, y_0) , i.e., replace each
- 234 image $I \in \mathcal{S}$ by a shifted image $I' = T_{x_0, y_0}(I)$ for which $I'(x, y) = I(x - x_0, y - y_0)$,
- then, we should get, in effect, the exact same set of images:

$$T_{x_0, y_0}(\mathcal{S}) \stackrel{\text{def}}{=} \{T(x_0, y_0)(I) : I \in \mathcal{S}\} \approx \mathcal{S}. \quad (24)$$

The only difference between these two sets of images may be on the few images where
 the cat is right at the image's boundary; in this paper, we will ignore this difference
 – just like we ignored the bounded-ness in the previous text. In this ignoring-bounds
 approximation, we conclude that

$$T_{x_0, y_0}(\mathcal{S}) = \{T(x_0, y_0)(I) : I \in \mathcal{S}\} = \mathcal{S}. \quad (25)$$

235 How does shift of the original image affect the input signals to the following
 236 convolution layers? In between the very first input layer and the following convolution
 237 layers, we may have (and usually do have) layers that perform “compression” of the
 238 (x, y) part – i.e., that transform:

- 239 • values corresponding to *several* points (x, y)
- 240 • into values corresponding to a *single* new point (x', y') .

241 In general, the (x, y) -shift of the original data corresponds to a shift of the transformed
 242 data – but by smaller shift values. For example, if data corresponding to each new
 243 (x, y) -point comes from data from four different “pre-compression” points, then the shift
 244 by (x_0, y_0) in the pre- (x, y) -compression layer corresponds to a shift of the convolution
 245 layer input by $(x_0/2, y_0/2)$.

Since the set of input images should not change if we apply a shift, we can conclude
 that for each convolution layer, the set of the corresponding inputs to this layer should
 also not change if we shift all these inputs, i.e., if we replace each input $F_d(x, y)$ with a
 shifted input

$$F'_d(x, y) \stackrel{\text{def}}{=} F_d(x - x_0, y - y_0) \quad (26)$$

246 for some shift (x_0, y_0) .

The set of inputs on which we compare different methods does not change when
 we apply a shift. So, if one method was better when we processed original inputs, it
 should still be better if we process shifted inputs – since the resulting set of inputs is the
 same. In other words, the quality (e.g., accuracy) $Q_{\mathcal{F}}(\mathcal{S})$ of a method corresponding
 to the family \mathcal{F} , when gauged by the set of inputs corresponding to original images
 should be the same as this method's quality $Q_{\mathcal{F}}(T_{x_0, y_0}(\mathcal{S}))$ on the set

$$T_{x_0, y_0}(\mathcal{S}) = \{T_{x_0, y_0}(F_d) : F_d \in \mathcal{S}\} \quad (27)$$

247 of all the inputs obtained from the original set \mathcal{S} by this shift – since these two sets of
 248 inputs are, in effect, the same set: $T_{x_0, y_0}(\mathcal{S}) = \mathcal{S}$. Thus, $Q_{\mathcal{F}}(T_{x_0, y_0}(\mathcal{S})) = Q_{\mathcal{F}}(\mathcal{S})$.

But, as one can see, shifting all the inputs is equivalent to shifting all the sets from the family \mathcal{F} . Indeed, if we apply the formula (23) to the shifted layer's input $F'_d(x, y) \stackrel{\text{def}}{=} F_d(x - x_0, y - y_0)$, we get

$$G_d(x, y) = \sum_{d'=1}^{d_{\text{in}}} \left(\sum_{(x', y') \in S_{\mathcal{F}}(x, y): d_{\infty}((x, y), (x', y')) \leq L} K_d(x, x', y, y', d') \cdot F_{d'}(x' - x_0, y' - y_0) \right), \quad (28)$$

i.e., in terms of the shifted coordinates $X \stackrel{\text{def}}{=} x - x_0$ and $Y \stackrel{\text{def}}{=} y - y_0$ for which $x = X + x_0$ and $y = Y + y_0$, we get – taking into account that the distance d_{∞} does change with shift – that:

$$G_d(X, Y) = \sum_{d'=1}^{d_{\text{in}}} \left(\sum_{C': d_{\infty}((X, Y), (X', Y')) \leq L} K'_d(X, X', Y, Y', d') \cdot F_{d'}(x' - x_0, y' - y_0) \right), \quad (29)$$

where we denoted

$$K'_d(X, X', Y, Y', d') \stackrel{\text{def}}{=} K_d(X + x_0, X' + x_0, Y + y_0, Y' + y_0, d'), \quad (30)$$

and where C' denotes the condition $(X' + x_0, Y' + y_0) \in S_{\mathcal{F}}(X + x_0, Y + y_0)$.

In terms of the family \mathcal{F} , the main difference between the formulas (23) and (29) is that instead of the condition $(x', y') \in S_{\mathcal{F}}(x, y)$, we now have a new condition

$$C' \Leftrightarrow (X' + x_0, Y' + y_0) \in S_{\mathcal{F}}(X + x_0, Y + y_0), \quad (31)$$

i.e., equivalently, $(X', Y') \in S_{\mathcal{F}}(X + x_0, Y + y_0) - (x_0, y_0)$. It is easy to check that this new condition is equivalent to $(Y', Y') \in S_{\mathcal{F}'}(X, Y)$, where the new family \mathcal{F}' is obtained by shifting sets from the original family \mathcal{F} .

So:

- the relative quality of two families does not change if we shift all the layer's inputs;
- however, shifting all the layer's inputs is equivalent to shifting all the sets from the family \mathcal{F} .

Thus, the relative quality of two families does not change if we shift both families. In other words, a reasonable optimality criterion – that describes which family is better – should be invariant with respect to shifts.

Similarly, we can argue that a reasonable optimality criterion should not change if we rename x - and y -axes, etc.

We are ready. Now, we are ready for the precise formulation of the problem.

3. Definitions and the Main Result

Definition 1. By a family, we mean a family of non-empty subsets of the “grid” $\mathbb{Z} \times \mathbb{Z}$, a family in which:

- all sets from this family are disjoint, and
- at least one set from this family has more than one element.

Terminological comment. To avoid possible misunderstandings, let us emphasize that here, we consider several levels of sets, and to avoid confusion, we use different terms for sets from different levels:

- first, we consider *points* $(x, y) \in \mathbb{Z} \times \mathbb{Z}$;
- second, we consider *sets* of points $S \subseteq \mathbb{Z} \times \mathbb{Z}$; we call them simply *sets*;
- third, we consider sets of sets of points $\mathcal{F} = \{S, S', \dots\}$; we call them *families*;
- finally, we consider the set of all possible families $\{\mathcal{F}, \mathcal{F}', \dots\}$; we call this a *class*.

275 *Comment about the requirements.* In the previous text, we argued that for each family
 276 \mathcal{F} , the union of all its sets $\cup\{S : S \in \mathcal{F}\}$ should coincide with the whole “grid” $\mathbb{Z} \times \mathbb{Z}$.
 277 However, in our definition of an alternative, we did not impose this requirement. We
 278 omitted this requirement to make our result stronger – since, as we see from the following
 279 Proposition, this requirement actually follows from all the other requirements.

280 **Definition 2.** *By an optimality criterion, we mean a pair of relations $(<, \sim)$ on the class of all*
 281 *possible families that satisfy the following conditions:*

- 282 • if $\mathcal{F} < \mathcal{F}'$ and $\mathcal{F}' < \mathcal{F}''$, then $\mathcal{F} < \mathcal{F}''$;
- 283 • if $\mathcal{F} < \mathcal{F}'$ and $\mathcal{F}' \sim \mathcal{F}''$, then $\mathcal{F} < \mathcal{F}''$;
- 284 • if $\mathcal{F} \sim \mathcal{F}'$ and $\mathcal{F}' < \mathcal{F}''$, then $\mathcal{F} < \mathcal{F}''$;
- 285 • if $\mathcal{F} \sim \mathcal{F}'$ and $\mathcal{F}' \sim \mathcal{F}''$, then $\mathcal{F}' \sim \mathcal{F}''$;
- 286 • we have $\mathcal{F} \sim \mathcal{F}$ for all \mathcal{F} ; and
- 287 • if $\mathcal{F} < \mathcal{F}'$, then we cannot have $\mathcal{F} \sim \mathcal{F}'$.

288 *Comment.* The pair of relations $(<, \sim)$ between families of subsets forms what is called a
 289 *pre-order* or *quasi-order*. This notion is more general than partial order, since, in contrast
 290 to the definition of the partial order, we do not require that if $a \leq b$ and $b \leq a$, then $a = b$:
 291 in principle, we can have $a \sim b$ for some $a \neq b$.

292 **Definition 3.** *We say that a family \mathcal{F} is optimal with respect to the optimality criterion $(<, \sim)$*
 293 *if for every other family \mathcal{F}' , we have either $\mathcal{F}' < \mathcal{F}$ or $\mathcal{F}' \sim \mathcal{F}$.*

294 **Definition 4.** *We say that the optimality criterion is final if there exists exactly one family*
 295 *which is optimal with respect to this criterion.*

296 **Definition 5.** *By a transformation $T : \mathbb{Z} \times \mathbb{Z}$, we mean one of the following transformations:*
 297 $T_{x_0, y_0}(x, y) = (x - x_0, y - y_0)$, $T_{-+}(x, y) = (-x, y)$, $T_{+-}(x, y) = (x, -y)$, and $T_{\leftrightarrow}(x, y) =$
 298 (y, x) .

299 **Definition 6.** *For each family \mathcal{F} and for each transformation T , by the result $T(\mathcal{F})$ of applying*
 300 *the transformation T to the family \mathcal{F} , we mean the family $T(\mathcal{F}) = \{T(S) : S \in \mathcal{F}\}$, where, for*
 301 *any set S , $T(S) \stackrel{\text{def}}{=} \{T(x, y) : (x, y) \in S\}$.*

302 **Definition 7.** *We say that the optimality criterion is invariant if for all transformations T ,*
 303 *$\mathcal{F} < \mathcal{F}'$ implies that $T(\mathcal{F}) < T(\mathcal{F}')$, and $\mathcal{F} \sim \mathcal{F}'$ implies that $T(\mathcal{F}) \sim T(\mathcal{F}')$.*

304 **Proposition.** *For every final invariant optimality criterion, the optimal family is equal, for some*
 305 *integer $\ell \geq 1$, to one of the following two families:*

- 306 • the family of all the sets $S_{\ell, x_0, y_0} \stackrel{\text{def}}{=} \{(x_0 + \ell \cdot n_x, y_0 + \ell \cdot n_y) : n_x, n_y \in \mathbb{Z}\}$ corresponding
 307 to all possible pairs of integers (x_0, y_0) for which $0 \leq x_0, y_0 < \ell$;
- the family of all the sets

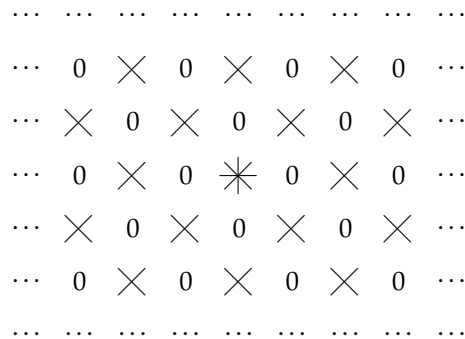
$$S'_{\ell, x_0, y_0} \stackrel{\text{def}}{=} \{(x_0 + \ell \cdot n_x, y_0 + \ell \cdot n_y) : n_x, n_y \in \mathbb{Z} \text{ and } n_x + n_y \text{ is even}\}$$

308 corresponding to all possible pairs of integers (x_0, y_0) for which $0 \leq x_0, y_0 < \ell$.

309 *Comments.*

- 310 • This proposition takes care of all invariant (and final) optimality criteria. Thus, it
 311 should work for all usual criteria based on misclassification rate, time of calculation,
 312 used memory, or any other used in neural networks: indeed, if one method is better
 313 than another for images in general, it should remain to be better if we simply shift
 314 all the images or turn all the images upside down. Images can come as they are, they
 315 can come upside down, they can come shifted, etc. If for some averaging criterion,
 316 one method works better for all possible images but another method works better
 317 for all upside-down versions of these images – which is, in effect, the same class of
 318 possible images – then from the common sense viewpoint, this would mean that
 319 something is not right with this criterion.

- 320 • The first possibly optimal case corresponds to dilated convolution. In the second
 321 possibly optimal case, the optimal family contains similar but somewhat different
 322 sets; an example of such a set is given in Fig. 7.



323
324
325
326 Figure 7. A set from the second possibly optimal family

327 Thus, this result explains the effectiveness of dilated convolution – and also provides
 328 us with a new alternative worth trying.

- 329 • The following proof is similar to several proofs presented in [4].

330 **Proof.**

331 1°. Since the optimality criterion is final, there exists exactly one optimal family \mathcal{F}_{opt} .
 332 Let us first prove that this family is itself invariant, i.e., that $T(\mathcal{F}_{\text{opt}}) = \mathcal{F}_{\text{opt}}$ for all
 333 transformations T .

334 Indeed, the fact that the family \mathcal{F}_{opt} is optimal means that for every family \mathcal{F} , we
 335 have $\mathcal{F} < \mathcal{F}_{\text{opt}}$ or $\mathcal{F} \sim \mathcal{F}_{\text{opt}}$. Since this is true for every family \mathcal{F} , it is also true for
 336 every family $T^{-1}(\mathcal{F})$, where T^{-1} denotes inverse transformation (i.e., a transformation
 337 for which $T(T^{-1}(x, y)) = (x, y)$). Thus, for every family \mathcal{F} , we have either $T^{-1}(\mathcal{F}) <$
 338 \mathcal{F}_{opt} or $T^{-1}(\mathcal{F}) \sim \mathcal{F}_{\text{opt}}$. Due to invariance, we have $\mathcal{F} = T(T^{-1}(\mathcal{F})) < T(\mathcal{F}_{\text{opt}})$ or
 339 $\mathcal{F} \sim T(\mathcal{F}_{\text{opt}})$. By definition of optimality, this means that the alternative $T(\mathcal{F}_{\text{opt}})$ is also
 340 optimal. However, since the optimality criterion is final, there exists exactly one optimal
 341 family, so $T(\mathcal{F}_{\text{opt}}) = \mathcal{F}_{\text{opt}}$.

342 The statement is proven.

343 2°. Let us now prove that the optimal family contains a set S' that, in its turn, contains
 344 the point $(0, 0)$ and some point $(x, y) \neq (0, 0)$.

Indeed, by definition of a family, every family – including the optimal family –
 contains at least one set S that has at least two points. Let S be any such set from the
 optimal family, and let (x_0, y_0) be any of its points. Then, due to Part 1 of this proof, the
 set $S' \stackrel{\text{def}}{=} T_{x_0, y_0}(S)$ also belongs to the optimal family, and this set contains the point

$$T_{x_0, y_0}(x_0, y_0) = (x_0 - x_0, y_0 - y_0) = (0, 0).$$

345 Since the set S had at least two different points, the set $S' = T_{x_0, y_0}(S)$ also contains
 346 at least two different points. Thus, the set S' must contain a point (x, y) which is different
 347 from $(0, 0)$.

348 The statement is proven.

349 3°. In the following text, by S' , we will mean a set from the optimal family \mathcal{F}_{opt} whose
 350 existence is proven in Part 2 of this proof: namely, a set that contains the point $(0, 0)$ and
 351 a point $(x, y) \neq (0, 0)$.

352 4°. Let us prove that if the set S' contains a point (x, y) , then this set also contains the
 353 points $(x, -y)$, $(-x, y)$, and (y, x) .

354 Indeed, due to Part 1 of this proof, with the set S' the optimal family \mathcal{F}_{opt} also
 355 contains the set $T_{+-}(S')$. This set contains the point $T_{+-}(0, 0) = (0, 0)$. Thus, the sets S'

356 and $T_{+-}(S')$ have a common element $(0,0)$. Since different sets from the optimal family
 357 must be disjoint, it follows that the sets S' and $T_{+-}(S')$ must coincide. The set $T_{+-}(S')$
 358 contains the points $(x, -y)$ for each point $(x, y) \in S$. Since $T_{+-}(S') = S'$, this implies
 359 that for each point $(x, y) \in S'$, we have $(x, -y) \in T_{+-}(S') = S'$.

360 Similarly, we can prove that $(-x, y) \in S'$ and $(y, x) \in S'$. The statement is proven.

5°. By combining the two conclusions of Part 4 – that $(x, -y) \in S'$ and that therefore
 $T_{-+}(x, -y) = (-x, -y) \in S'$, we conclude that for every point $(x, y) \in S'$, the point

$$-(x, y) \stackrel{\text{def}}{=} (-x, -y)$$

361 is also contained in the set S' .

6°. Let us prove that if the set S' contains two points (x_1, y_1) and (x_2, y_2) , then it also
 contains the point

$$(x_1, y_1) + (x_2, y_2) \stackrel{\text{def}}{=} (x_1 + x_2, y_1 + y_2).$$

Indeed, due to Part 1 of this proof, the set $T_{-x_2, -y_2}(S')$ also belongs to the optimal
 family. This set shares an element

$$T_{-x_2, -y_2}(0,0) = (0 - (-x_2), 0 - (-y_2)) = (x_2, y_2)$$

with the original set S' . Thus, the set $T_{-x_2, -y_2}(S')$ must coincide with the set S' . Due to
 the fact that $(x_1, y_1) \in S'$, the element

$$T_{-x_2, -y_2}(x_1, y_1) = (x_1 - (-x_2), y_1 - (-y_2)) = (x_1 + x_2, y_1 + y_2)$$

362 belongs to the set $T_{x_1, y_1}(S') = S'$. The statement is proven.

7°. Let us prove that if the set S' contains a point (x, y) , then, for each integer c , this set
 also contains the point

$$c \cdot (x, y) = (c \cdot x, c \cdot y).$$

Indeed, if c is positive, this follows from the fact that

$$(c \cdot x, c \cdot y) = (x, y) + \dots + (x, y) \text{ (} c \text{ times)}.$$

363 When c is negative, then we first use Part 5 and conclude that $(-x, -y) \in S'$, and then
 364 conclude that the point $(|c| \cdot (-x), |c| \cdot (-y)) = (c \cdot x, c \cdot y)$ is in the set S' .

8°. Let us prove that if the set S' contains points $(x_1, y_1), \dots, (x_n, y_n)$, then for all integers
 c_1, \dots, c_n , it also contains their linear combination

$$c_1 \cdot (x_1, y_1) + \dots + c_n \cdot (x_n, y_n) = (c_1 \cdot x_1 + \dots + c_n \cdot x_n, c_1 \cdot y_1 + \dots + c_n \cdot y_n).$$

365 Indeed, this follows from Parts 6 and 7.

366 9°. The set S' contains some points which are different from $(0,0)$, i.e., points for which
 367 at least one of the integer coordinates is non-zero. According to Parts 4 and 5, we can
 368 change the signs of both x and y coordinates and still get points from S' . Thus, we can
 369 always consider points with non-negative coordinates.

370 Let d denote the greatest common divisor of all positive values of the coordinates
 371 of points from S' .

If a value x appears as an x -coordinate of some point $(x, y) \in S'$, then, due to Part 4,
 we have $(x, -y) \in S'$ and thus, due to Part 5,

$$(x, y) + (x, -y) = (2x, 0) \in S'.$$

372 Similarly, if a value y appears as a y -coordinate of some point $(x, y) \in S'$, then we get
 373 $(0, 2y) \in S'$ and thus, due to Part 3, $(2y, 0) \in S'$.

It is known that a common divisor d of the values v_1, \dots, v_n can be represented as a linear combination of these values:

$$d = c_1 \cdot v_1 + \dots + c_n \cdot v_n.$$

For each value v_i , we have $(2v_i, 0) \in S'$, thus, for

$$2d = c_1 \cdot (2v_1) + \dots + c_n \cdot (2v_n),$$

374 by Part 8, we get $(2d, 0) \in S'$. Due to Part 4, we thus have $(0, 2d) \in S'$, and due to Parts
375 6 and 7, all points $(n_x \cdot (2d), n_y \cdot (2d))$ for integers n_x and n_y also belong to the set S' .

376 If S' has no other points, then for the set containing $(0, 0)$, we indeed conclude that
377 this set is the same as what we described for dilated convolution, with $\ell = 2d$.

378 10°. What if these are other points in the set S' ? Since d is the greatest common divisor
379 of all the coordinate values, each of these points has the form $(c_x \cdot d, c_y \cdot d)$ for some
380 integers c_x and c_y . Since this point is not of the form $(n_x \cdot (2d), n_y \cdot (2d))$, this means that
381 either c_x , or c_y is an odd number – or both.

Let us first consider the case when exactly one of the values c_x and c_y is odd. Without losing generality, let us assume that c_x is odd, so $c_x = 2n_x + 1$ and $c_y = 2n_y$ for some integers n_x and n_y . Due to Part 9, we have $(2n_x \cdot d, 2n_y \cdot d) \in S'$, so the difference

$$((2n_x + 1) \cdot d, 2n_y \cdot d) - (2n_x \cdot d, 2n_y \cdot d) = (d, 0)$$

382 also belongs to the set S' . Thus, similarly to Part 9, we can conclude that for every two
383 integers c_x and c_y , we have $(c_x \cdot d, c_y \cdot d) \in S'$. So, in this case, S' coincides, for $\ell = d$,
384 with the set corresponding to dilated convolution.

The only remaining case is when not all points $(c_x \cdot d, c_y \cdot d)$ belong to the set S' . This means that for some such point both values c_x and c_y are odd: $c_x = 2n_x + 1$ and $c_y = 2n_y + 1$ for some integers n_x and n_y . Due to Part 9, we have $(2n_x \cdot d, 2n_y \cdot d) \in S'$, so the difference

$$((2n_x + 1) \cdot d, (2n_y + 1) \cdot d) - (2n_x \cdot d, 2n_y \cdot d) = (d, d)$$

385 also belongs to the set S' .

386 Since, due to Part 9, we have $(2n_x \cdot d, 2n_y \cdot d) \in S'$ for all n_x and n_y , we conclude,
387 by using Part 5, that $((2n_x + 1) \cdot d, (2n_y + 1) \cdot d) \in S'$. So, all pairs for which both
388 coordinates are odd multiples of d are in S' . Thus, we get the new case described in the
389 Proposition.

390 11°. The previous results were about the sets containing the point $(0, 0)$.

391 For all other sets S containing some other point (x_0, y_0) , we get the same result if
392 we take into account that the optimal family is invariant, and thus, with the set S , the
393 optimal family also contains the set $T_{x_0, y_0}(S)$ that contains $(0, 0)$ and is, thus, equal either
394 to the family corresponding to dilated convolution or to the new similar family.

395 The proposition is proven.

396 4. Conclusions and Future Work

397 **Conclusions.** One of the efficient machine learning ideas is the idea of a convolutional
398 neural network. Such networks use convolutional layers, in which the layer's output at
399 each point is a combination of the layer's input corresponding to several neighboring
400 points. A reasonable idea is to restrict ourselves to only some of the neighboring points.
401 It turns out that out of all such restrictions, the best results are obtained when we only use
402 neighboring points for which the differences in both coordinates are divisible by some
403 constant ℓ . Networks implementing such restrictions are known as dilated convolutional
404 neural networks.

405 In this paper, we provide a theoretical explanation for this empirical conclusion.

4.06 **Future work.** This paper describes a general abstract result: that for any optimality
4.07 criterion that satisfies some reasonable properties, *some* dilated convolution is optimal.
4.08 To be practically useful, it is desirable to analyze which dilated convolutions are optimal
4.09 for different practical situations and for specific criteria uses in machine learning, such
4.10 as misclassification rate, time of calculation, used memory, etc. (nd the combination of
4.11 these criteria). It is also desirable to analyze what size neighborhood should we choose
4.12 for different practical situations and for different criteria.

4.13 **Author Contributions:** All three authors contributed equally to this paper. All three authors have
4.14 read and agreed to the published version of the manuscript.

4.15 **Funding:** This work was supported in part by the National Science Foundation grants 1623190 (A
4.16 Model of Change for Preparing a New Generation for Professional Practice in Computer Science),
4.17 and HRD-1834620 and HRD-2034030 (CAHSI Includes). It was also supported by the program of
4.18 the development of the Scientific-Educational Mathematical Center of Volga Federal District No.
4.19 075-02-2020-1478.

4.20 **Acknowledgments:** The authors are greatly thankful to the anonymous referees for valuable
4.21 suggestions.

4.22 **Conflicts of Interest:** The authors declare no conflict of interest.

References

1. Goodfellow, I.; Bengio, Y.; Courville, A. *Deep Learning*, MIT Press: Cambridge, Massachusetts, 2016.
2. Kreinovich, V.; Kosheleva, O. Optimization under uncertainty explains empirical success of deep learning heuristics", In: Pardalos, P.; Rasskazova, V.; Vrahatis, M.N. (eds.), *Black Box Optimization, Machine Learning and No-Free Lunch Theorems*, Springer: Cham, Switzerland, 2021, pp. 195–220.
3. Li, Y.; Zhang, X.; Chen, D. CSRNet: dilated convolutional neural networks for understanding the highly congested scenes, *Proceedings of the 2018 Conference on Computer Vision and Pattern Recognition CVPR'2018*, Salt Lake City, Utah, June 18–22, 2018, pp. 1091–1100.
4. Nguyen, H.T.; Kreinovich, V. *Applications of Continuous Mathematics to Computer Science*, Kluwer: Dordrecht, 1997.
5. Yu, F.; Koltun, V. Multi-scale context aggregation by dilated convolutions, *Proceedings of the 4th International Conference on Learning Representations ICLR'2016*, San Juan, Puerto Rico, May 2–4, 2016.
6. Zhang, X.; Zou, Y.; Shi, W. Dilated convolution neural network with LeakyReLU for environmental sound classification, *Proceedings of the 2017 22nd International Conference on Digital Signal Processing DSP'2017*, London, U.K., August 23–25, 2017.