

1 **SUPPLEMENTARY MATERIALS: ABBA Neural Networks: Coping with Positivity, Expressivity,**
2 **and Robustness ***

3 Ana Neacșu [†], Jean-Christophe Pesquet [‡], Vlad Vasilescu [†], and Corneliu Burileanu [†]
4

5 **SM1. Symmetric activation functions.** In practice, the activation operator R_i is often separa-
6 ble, that is it operates componentwise:

7 (SM1.1)
$$(\forall x = (\xi_k)_{1 \leq k \leq N_i} \in \mathbb{R}^{N_i}) \quad R_i x = (\varrho_i(\xi_k))_{1 \leq k \leq N_i},$$

8 where, for every $k \in \{1, \dots, N_i\}$, $\varrho_i: \mathbb{R} \rightarrow \mathbb{R}$. Examples of odd functions allowing us to define
9 a symmetric separable activation operators R_i with $c_i = d_i = 0$ are

- 10 • the hyperbolic tangent activation function $\rho_i = \tanh$
11 • the arctangent activation function $\rho_i = (2/\pi) \arctan$
12 • the inverse square root linear unit function $\varrho_i: \mathbb{R} \rightarrow \mathbb{R}: \xi \mapsto \xi/\sqrt{1 + \xi^2}$
13 • the Elliot activation function $\varrho_i: \mathbb{R} \rightarrow \mathbb{R}: \xi \mapsto \xi/(1 + |\xi|)$.

14 Some examples of separable activation operators which are non-odd are described below. The
15 capped ReLU function is given by

16 (SM1.2)
$$(\forall \xi \in \mathbb{R}) \quad \rho_i(\xi) = \begin{cases} 0 & \text{if } \xi < 0 \\ \xi & \text{if } 0 \leq \xi < \chi \\ \chi & \text{otherwise,} \end{cases}$$

17 where $\chi \in]0, +\infty[$. We have then $c_i = d_i = \chi \mathbf{1}_{N_i}$ with $\mathbf{1}_{N_i} = [1, \dots, 1]^T \in \mathbb{R}^{N_i}$. We can also
18 define a leaky version of this function as

19 (SM1.3)
$$(\forall \xi \in \mathbb{R}) \quad \rho_i(\xi) = \begin{cases} \alpha \xi & \text{if } \xi < 0, \\ \xi & \text{if } 0 \leq \xi < \chi, \\ \alpha(\xi - \chi) + \chi & \text{otherwise,} \end{cases}$$

20 where $\alpha \in]0, 1[$ and $\chi \in]0, +\infty[$ are hyper-parameters.

21 **SM2. Proof of the properties of ABBA matrices.** (i)-(iii): These properties follow from basic
22 algebra. We will just detail the proof of the third one. Let

23 (SM2.1)
$$M_1 = \begin{bmatrix} A_1 & B_1 \\ B_1 & A_1 \end{bmatrix} \quad \text{and} \quad M_2 = \begin{bmatrix} A_2 & B_2 \\ B_2 & A_2 \end{bmatrix},$$

24 where $(A_1, B_1) \in \mathbb{R}^{N_2 \times N_1}$ and $(A_2, B_2) \in \mathbb{R}^{N_3 \times N_2}$. Then

25 (SM2.2)
$$M_2 M_1 = \begin{bmatrix} A_2 A_1 + B_2 B_1 & A_2 B_1 + B_2 A_1 \\ A_2 B_1 + B_2 A_1 & A_2 A_1 + B_2 B_1 \end{bmatrix} \in \mathcal{A}_{N_3, N_1}.$$

*Part of this work was supported by the French ANR Research and Teaching Chair in Artificial Intelligence BRIDGEABLE.

[†]Speech and Dialogue Laboratory, University Politehnica of Bucharest, Romania (ana_antonia.neacsu@upb.ro).

[‡]Centre de Vision Numérique, Inria, CentraleSupélec, Gif-sur-Yvette, France.

26 In addition,

$$\begin{aligned}
27 \quad & \mathfrak{S}(M_2 M_1) = A_2 A_1 + B_2 B_1 + A_2 B_1 + B_2 A_1 \\
28 \quad & = (A_2 + B_2)(A_1 + B_1) \\
30 \quad \text{(SM2.3)} \quad & = \mathfrak{S}(M_2) \mathfrak{S}(M_1).
\end{aligned}$$

31 (iv): This property is a direct consequence of (ii) and (iii).

32 (v): Let $M = \begin{bmatrix} A & B \\ B & A \end{bmatrix} \in \mathbb{R}^{(2N_1) \times (2N_1)}$. $\lambda \in \mathbb{C}$ is an eigenvalue of M if and only if

$$33 \quad \text{(SM2.4)} \quad \det(M - \lambda \text{Id}) = 0 \Leftrightarrow \det \left(\begin{bmatrix} A - \lambda \text{Id} & B \\ B & A - \lambda \text{Id} \end{bmatrix} \right) = 0.$$

34 We have

$$35 \quad \text{(SM2.5)} \quad \begin{bmatrix} A - \lambda \text{Id} & B \\ B & A - \lambda \text{Id} \end{bmatrix} \begin{bmatrix} \text{Id} & -\text{Id} \\ \text{Id} & \text{Id} \end{bmatrix} = \begin{bmatrix} A + B - \lambda \text{Id} & -A + B + \lambda \text{Id} \\ A + B - \lambda \text{Id} & A - B - \lambda \text{Id} \end{bmatrix}.$$

36 Since $A - B - \lambda \text{Id}$ and $-A + B + \lambda \text{Id}$ commute, we have [SM5]

$$37 \quad \text{(SM2.6)} \quad \det \left(\begin{bmatrix} A + B - \lambda \text{Id} & -A + B + \lambda \text{Id} \\ A + B - \lambda \text{Id} & A - B - \lambda \text{Id} \end{bmatrix} \right) = 2^N \det((A + B - \lambda \text{Id})(A - B - \lambda \text{Id})).$$

38 Similarly

$$39 \quad \text{(SM2.7)} \quad \det \left(\begin{bmatrix} \text{Id} & -\text{Id} \\ \text{Id} & \text{Id} \end{bmatrix} \right) = 2^N.$$

40 We deduce from (SM2.5) that

$$\begin{aligned}
41 \quad & \det \left(\begin{bmatrix} A - \lambda \text{Id} & B \\ B & A - \lambda \text{Id} \end{bmatrix} \right) = \det((A + B - \lambda \text{Id})(A - B - \lambda \text{Id})) \\
43 \quad \text{(SM2.8)} \quad & \Leftrightarrow \det(M - \lambda \text{Id}) = \det(A + B - \lambda \text{Id}) \det(A - B - \lambda \text{Id}).
\end{aligned}$$

44 So λ is an eigenvalue of M if and only if $\det(A + B - \lambda \text{Id}) = 0$ or $\det(A - B - \lambda \text{Id}) = 0$, i.e.,
45 λ is an eigenvalue of $A + B$ or $A - B$.

46 (vi) Let M be defined similarly to previously with $(A, B) \in (\mathbb{R}^{N_2 \times N_1})^2$. We have

$$47 \quad \text{(SM2.9)} \quad \|M\| = \|MM^\top\|^{1/2} = \left\| \begin{bmatrix} AA^\top + BB^\top & AB^\top + BA^\top \\ AB^\top + BA^\top & AA^\top + BB^\top \end{bmatrix} \right\|^{1/2}$$

48 According to (v), the eigenvalues of $MM^\top \in \mathcal{A}_{N_2, N_2}$ are those of $AA^\top + BB^\top + AB^\top + BA^\top =$
49 $(A + B)(A + B)^\top$ and $AA^\top + BB^\top - AB^\top - BA^\top = (A - B)(A - B)^\top$. The maximum
50 eigenvalues of the two latter matrices are $\|A + B\|^2$ and $\|A - B\|^2$, respectively. Therefore

$$51 \quad \|M\| = \max\{\|A + B\|, \|A - B\|\}.$$

52 (vii): In addition, if A and B have nonnegative elements,

$$\begin{aligned} 53 \quad \|A - B\| &= \sup_{x \in \mathbb{R}^N \setminus \{0\}} \frac{\|Ax - Bx\|}{\|x\|} \\ 54 &\leq \sup_{x \in \mathbb{R}^N \setminus \{0\}} \frac{\|A|x| + B|x|\|}{\|x\|} \\ 55 &= \sup_{a \in [0, +\infty[^N \setminus \{0\}} \frac{\|Aa + Ba\|}{\|a\|} \\ 56 \quad (\text{SM2.10}) &\leq \|A + B\|, \end{aligned}$$

58 where $|x|$ denotes the vector whose components are the absolute values of those of vector x .

59 We deduce from (vi) that $\|M\| = \|A + B\| = \|\mathfrak{S}(M)\|$.

60 (viii): We have

$$61 \quad (\text{SM2.11}) \quad A + B = \sum_{k=1}^K \lambda_k u_k v_k^\top$$

$$62 \quad (\text{SM2.12}) \quad A - B = \sum_{k=1}^K \mu_k t_k w_k^\top.$$

64 Thus

$$65 \quad (\text{SM2.13}) \quad A = \frac{1}{2} \sum_{k=1}^K (\lambda_k u_k v_k^\top + \mu_k t_k w_k^\top)$$

$$66 \quad (\text{SM2.14}) \quad B = \frac{1}{2} \sum_{k=1}^K (\lambda_k u_k v_k^\top - \mu_k t_k w_k^\top)$$

68 and we deduce that

$$\begin{aligned} 69 \quad \begin{bmatrix} A & B \\ B & A \end{bmatrix} &= \sum_{k=1}^K \frac{1}{2} \left(\lambda_k \begin{bmatrix} u_k v_k^\top & u_k v_k^\top \\ u_k v_k^\top & u_k v_k^\top \end{bmatrix} + \mu_k \begin{bmatrix} t_k w_k^\top & -t_k w_k^\top \\ -t_k w_k^\top & t_k w_k^\top \end{bmatrix} \right) \\ 70 \quad (\text{SM2.15}) &= \sum_{k=1}^K \frac{1}{2} \left(\lambda_k \begin{bmatrix} u_k \\ u_k \end{bmatrix} \begin{bmatrix} v_k \\ v_k \end{bmatrix}^\top + \mu_k \begin{bmatrix} t_k \\ -t_k \end{bmatrix} \begin{bmatrix} w_k \\ -w_k \end{bmatrix}^\top \right). \\ 71 \end{aligned}$$

72 On the other hand, for every $(k, \ell) \in \{1, \dots, K\}^2$,

$$73 \quad \begin{bmatrix} u_k \\ u_k \end{bmatrix}^\top \begin{bmatrix} u_\ell \\ u_\ell \end{bmatrix} = 2u_k^\top u_\ell = \begin{cases} 2 & \text{if } k = \ell \\ 0 & \text{otherwise,} \end{cases}$$

$$74 \quad \begin{bmatrix} t_k \\ -t_k \end{bmatrix}^\top \begin{bmatrix} t_\ell \\ -t_\ell \end{bmatrix} = 2t_k^\top t_\ell = \begin{cases} 2 & \text{if } k = \ell \\ 0 & \text{otherwise,} \end{cases}$$

$$75 \quad (\text{SM2.16}) \quad \begin{bmatrix} u_k \\ u_k \end{bmatrix}^\top \begin{bmatrix} t_\ell \\ -t_\ell \end{bmatrix} = 0,$$

76

77 which shows that $\left\{ \frac{1}{\sqrt{2}} \begin{bmatrix} u_k \\ u_k \end{bmatrix}, \frac{1}{\sqrt{2}} \begin{bmatrix} t_k \\ -t_k \end{bmatrix} \right\}_{1 \leq k \leq K}$ is an orthonormal family of \mathbb{R}^{2N_2} . For similar
 78 reasons, $\left\{ \frac{1}{\sqrt{2}} \begin{bmatrix} v_k \\ v_k \end{bmatrix}, \frac{1}{\sqrt{2}} \begin{bmatrix} w_k \\ -w_k \end{bmatrix} \right\}_{1 \leq k \leq K}$ is an orthonormal family of \mathbb{R}^{2N_1} . This allows us to
 79 conclude that (SM2.15) provides a singular value decomposition of $\begin{bmatrix} A & B \\ B & A \end{bmatrix}$.

80 (ix): The rank of $\begin{bmatrix} A & B \\ B & A \end{bmatrix}$ is equal to the number of its nonzero singular values. From the
 81 previous result, it is thus equal to the sum of the nonzero values of $A + B$ and those of $A - B$,
 82 that is the sum of the ranks of matrices $A + B$ and $A - B$.

83 (x): The fact that the ABBA structure is kept by matrix mappings operating elementwise is
 84 obvious. Let us thus focus on the case of spectral functions. By using the same notation as in
 85 (viii), it follows from (SM2.15) that

(SM2.17)

$$86 \quad f\left(\begin{bmatrix} A & B \\ B & A \end{bmatrix}\right) = \sum_{k=1}^K \frac{1}{2} \left(\varphi(\lambda_k) \begin{bmatrix} u_k v_k^\top & u_k v_k^\top \\ u_k v_k^\top & u_k v_k^\top \end{bmatrix} + \varphi(\mu_k) \begin{bmatrix} t_k w_k^\top & -t_k w_k^\top \\ -t_k w_k^\top & t_k w_k^\top \end{bmatrix} \right) = \begin{bmatrix} \tilde{A} & \tilde{B} \\ \tilde{B} & \tilde{A} \end{bmatrix},$$

88 where

$$89 \quad \tilde{A} + \tilde{B} = \sum_{k=1}^K \varphi(\lambda_k) u_k v_k^\top$$

$$90 \quad \text{(SM2.18)} \quad \tilde{A} - \tilde{B} = \sum_{k=1}^K \varphi(\mu_k) t_k w_k^\top.$$

92 (xi): By using the same notation as in (3.6), The best approximation of rank less than or equal
 93 to R to a matrix M_0 in $\mathbb{R}^{(2N_2) \times (2N_1)}$ is $f(M_0)$ where f is given by (3.6) with

$$94 \quad \text{(SM2.19)} \quad (\forall \lambda \in \mathbb{R}_+) \quad \varphi(\lambda) = \begin{cases} \lambda & \text{if } \lambda \leq \tilde{\lambda}_{0,[R]} \\ 0 & \text{otherwise,} \end{cases}$$

95 and $\tilde{\lambda}_{0,[R]}$ is the R -th eigenvalue of M_0 when these are ordered by decreasing value: $\tilde{\lambda}_{0,1} \geq$
 96 $\dots \geq \tilde{\lambda}_{0,K}$. It thus follows from (x) that if $M_0 \in \mathcal{A}_{N_2, N_1}$, then $f(M_0) \in \mathcal{A}_{N_2, N_1}$.

97 (xii): The projection onto the spectral ball of center 0 and and radius $\rho \in]0, +\infty[$ of a matrix
 98 $M \in \mathcal{A}_{N_2, N_1}$ is given by (3.6) where

$$99 \quad (\forall \xi \in \mathbb{R}) \quad \varphi(\xi) = \min\{\xi, \rho\}.$$

100 The result then follows from Property (x).

101 *Remark SM2.1.* The last result can be generalized as follows. Let $\psi: \mathbb{R} \rightarrow]-\infty, +\infty]$ be a
 102 lower-semicontinuous function, which is proper, even, and convex, and let

$$103 \quad g: \mathbb{R}^{(2N_2) \times (2N_1)} \rightarrow]-\infty, +\infty]$$

$$104 \quad \text{(SM2.20)} \quad M \mapsto \sum_{i=1}^{2K} \psi(\tilde{\lambda}_k)$$

105

106 where $K = \min\{N_1, N_2\}$ and $(\tilde{\lambda}_k)_{1 \leq k \leq 2K}$ are the singular values of M . The proximity operator
 107 of g at $M \in \mathbb{R}^{(2N_2) \times (2N_1)}$ is [SM1, Proposition 24.68]:

$$\begin{aligned}
 108 \quad \text{prox}_g: M &\mapsto \underset{P \in \mathbb{R}^{(2N_2) \times (2N_1)}}{\text{argmin}} \frac{1}{2} \|P - M\|_F^2 + g(P) \\
 109 \quad (\text{SM2.21}) \quad &= \sum_{k=1}^{2K} \text{prox}_{\psi}(\tilde{\lambda}_k) \tilde{u}_k \tilde{v}_k^\top, \\
 110
 \end{aligned}$$

111 where $\|\cdot\|_F$ denotes the Frobenius norm. It then follows from Property (x) that, if $M \in \mathcal{A}_{N_2, N_1}$,
 112 then $\text{prox}_g(M) \in \mathcal{A}_{N_2, N_1}$.

113 **SM3. Link between Conv layers and MIMO systems.** To be rigorous, let us first define the
 114 space \mathcal{H}_{i-1} (resp. \mathcal{H}_i) in which signals $(x_p)_{1 \leq p \leq \zeta_{i-1}}$ (resp. $(y_q)_{1 \leq q \leq \zeta_i}$) used in (4.1) live.
 115 Typically, \mathcal{H}_i is some finite-dimensional subspace of $(\ell^2(\mathbb{Z}^d))^{\zeta_i}$ where $\ell^2(\mathbb{Z}^d)$ denotes the space
 116 of square summable discrete d -dimensional fields. For the discrete convolution $*$ to be properly
 117 defined, kernels $(w_{i,q,p})_{1 \leq p \leq \zeta_{i-1}, 1 \leq q \leq \zeta_i}$ are then assumed to be summable. In practice, this
 118 assumption is satisfied since these kernels are chosen with finite size.

119 For $x = (x(\mathbf{n}))_{\mathbf{n} \in \mathbb{Z}^d} \in \ell^2(\mathbb{Z}^d)$, the decimation operation $(\cdot) \downarrow_{s_i}$ returns the output signal

$$120 \quad (\text{SM3.1}) \quad (\forall \mathbf{n} \in \mathbb{Z}^d) \quad y(\mathbf{n}) = u(s_i \mathbf{n}).$$

121 Eq. (4.1) defines a MIMO (multi-input multi-output) filter that can be reexpressed in a matrix
 122 form as

$$\begin{aligned}
 123 \quad (\forall \mathbf{n} \in \mathbb{Z}^d) \quad \mathbf{u}(\mathbf{n}) &= \sum_{\mathbf{n}' \in \mathbb{Z}^d} \mathbf{W}_i(\mathbf{n}') \mathbf{x}(\mathbf{n} - \mathbf{n}') \\
 124 \quad (\text{SM3.2}) \quad &= (\mathbf{W}_i * \mathbf{x})(\mathbf{n}),
 \end{aligned}$$

126 where

$$127 \quad (\text{SM3.3}) \quad \mathbf{u}(\mathbf{n}) = \begin{bmatrix} u_1(\mathbf{n}) \\ \vdots \\ u_{\zeta_i}(\mathbf{n}) \end{bmatrix} \in \mathbb{R}^{\zeta_i}, \quad \mathbf{x}(\mathbf{n}) = \begin{bmatrix} x_1(\mathbf{n}) \\ \vdots \\ x_{\zeta_{i-1}}(\mathbf{n}) \end{bmatrix} \in \mathbb{R}^{\zeta_{i-1}},$$

129 and $\mathbf{W}_i(\mathbf{n})$ is given by (4.4). $(\mathbf{W}_i(\mathbf{n}))_{\mathbf{n} \in \mathbb{Z}^d}$ defines the so-called MIMO impulse response of
 130 \mathcal{W}_i . The MIMO impulse response of an ABBA layer is similarly given by (4.5).

131 These relations can also be written more concisely in the d -dimensional frequency domain¹
 132 as

$$133 \quad (\text{SM3.4}) \quad (\forall \boldsymbol{\nu} \in [0, 1]^d) \quad \hat{\mathbf{u}}(\boldsymbol{\nu}) = \widehat{\mathbf{W}}_i(\boldsymbol{\nu}) \hat{\mathbf{x}}(\boldsymbol{\nu}),$$

134 where

$$135 \quad (\text{SM3.5}) \quad \hat{\mathbf{x}}(\boldsymbol{\nu}) = \sum_{\mathbf{n} \in \mathbb{Z}^d} \mathbf{x}(\mathbf{n}) \exp(-i2\pi \mathbf{n}^\top \boldsymbol{\nu}) \in \mathbb{C}^{\zeta_{i-1}},$$

¹Alternatively, we could use the d -dimensional z-transform since we are dealing with discrete-space signals.

136

$$137 \quad (\text{SM3.6}) \quad \widehat{\mathbf{W}}_i(\boldsymbol{\nu}) = \sum_{\mathbf{n} \in \mathbb{Z}^d} \mathbf{W}_i(\mathbf{n}) \exp(-i2\pi \mathbf{n}^\top \boldsymbol{\nu}) \in \mathbb{C}^{\zeta_{j-1} \times \zeta_j},$$

138 and $\widehat{\mathbf{W}}_i$ is the frequency response of the associated MIMO filter.

139 Note that $\int_{[0,1]^d} \|\widehat{\mathbf{x}}(\boldsymbol{\nu})\|^2 d\boldsymbol{\nu} < +\infty$, whereas $\widehat{\mathbf{W}}_i$ is a continuous (hence bounded) function
 140 on $[0, 1]^d$. Another useful result from sampling theory [SM6] is that the Fourier transform of
 141 $\mathbf{y} = (y_q)_{1 \leq q \leq \zeta_j}$ in (4.1) is deduced from the Fourier transform of \mathbf{u} by the relation

$$142 \quad (\text{SM3.7}) \quad (\forall \boldsymbol{\nu} \in [0, 1]^d) \quad \widehat{\mathbf{y}}(\boldsymbol{\nu}) = \frac{1}{s_i^d} \sum_{\mathbf{j} \in \mathbb{S}(s_i)} \widehat{\mathbf{u}}\left(\frac{\boldsymbol{\nu} + \mathbf{j}}{s_i}\right).$$

143 where

$$144 \quad (\text{SM3.8}) \quad (\forall \sigma \in \mathbb{N} \setminus \{0\}) \quad \mathbb{S}(\sigma) = \{0, \dots, \sigma - 1\}^d.$$

145 It is also worth noting that the interpolation by a factor s of \mathbf{y}

$$146 \quad (\text{SM3.9}) \quad \mathbf{v} = \mathbf{y}_{\uparrow s} \Leftrightarrow (\forall \mathbf{n} \in \mathbb{Z}^d) \quad \mathbf{v}(\mathbf{n}) = \begin{cases} \mathbf{y}\left(\frac{\mathbf{n}}{s}\right) & \text{if } \mathbf{n} \in s\mathbb{Z}^d \\ 0 & \text{otherwise,} \end{cases}$$

147 translates into

$$148 \quad (\text{SM3.10}) \quad (\forall \boldsymbol{\nu} \in [0, 1]) \quad \widehat{\mathbf{y}}_{\uparrow s}(\boldsymbol{\nu}) = \widehat{\mathbf{y}}(s\boldsymbol{\nu}),$$

149 in the frequency domain.

150 **SM4. Frequency expressions of Lipschitz bounds.** In this appendix, we establish frequency-
 151 based bounds of the Lipschitz constant of an m -layer convolutional neural network T .

152 Based on the MIMO concepts introduced in Appendix SM3, we will introduce the following
 153 global frequency response of the network:

$$154 \quad (\text{SM4.1}) \quad (\forall \boldsymbol{\nu} \in [0, 1]^d) \quad \widehat{\mathbf{W}}(\boldsymbol{\nu}) = \widehat{\mathbf{W}}_m(\sigma_{m-1}\boldsymbol{\nu}) \cdots \widehat{\mathbf{W}}_2(\sigma_1\boldsymbol{\nu}) \widehat{\mathbf{W}}_1(\boldsymbol{\nu}) \in \mathbb{C}^{\zeta_m \times \zeta_0},$$

155 where $\widehat{\mathbf{W}}_i$ is the frequency response associated to filter \mathbf{W}_i (see (SM3.6)).

156 We have then the following result providing a frequency formula for evaluating the Lipschitz
 157 constant of a convolutional network.

158 **Proposition SM4.1.** *The quantity*

$$159 \quad (\text{SM4.2}) \quad \theta_m = \frac{1}{\sigma_m^{d/2}} \sup_{\boldsymbol{\nu} \in [0, 1/\sigma_m]^d} \left\| \sum_{\mathbf{j} \in \mathbb{S}(\sigma_m)} \widehat{\mathbf{W}}\left(\boldsymbol{\nu} + \frac{\mathbf{j}}{\sigma_m}\right) \widehat{\mathbf{W}}\left(\boldsymbol{\nu} + \frac{\mathbf{j}}{\sigma_m}\right)^H \right\|^{1/2}.$$

160 provides a lower bound on the Lipschitz constant estimate of network T ². In addition, if for every
 161 $i \in \{1, \dots, m\}$, $p \in \{1, \dots, \zeta_{i-1}\}$, and $q \in \{1, \dots, \zeta_i\}$, $w_{i,q,p} = (w_{i,q,p}(\mathbf{n}))_{\mathbf{n} \in \mathbb{Z}^d}$ is a nonnegative
 162 kernel i.e.,

$$163 \quad (\text{SM4.3}) \quad (\forall \mathbf{n} \in \mathbb{Z}^d) \quad w_{i,p,q}(\mathbf{n}) \geq 0,$$

164 then θ_m is a Lipschitz constant of T .

165 *Proof.* In the considered case all activation operators are nonexpansive and they are
 166 assumed separable, except maybe at the last layer. Thus T is a special case of the networks
 167 investigated in [SM3, Section 5] for which a tight estimate of the Lipschitz constant was
 168 provided. It then follows from [SM3, Theorem 5.2] that a lower bound on this Lipschitz
 169 constant estimate is

$$170 \quad (\text{SM4.4}) \quad \theta_m = \|\mathcal{W}_m \circ \dots \circ \mathcal{W}_1\|.$$

171 In addition, under the additional assumption that all the kernels are nonnegative, T is an
 172 instance of the positively weighted networks investigated in [SM3, Section 5.3] and it follows
 173 from [SM3, Proposition 5.10] that θ_m is then a Lipschitz constant of T .

174 So the problem is to calculate the norm of the linear operator $\mathcal{W} = \mathcal{W}_m \circ \dots \circ \mathcal{W}_1$. Each
 175 operator \mathcal{W}_i with $i \in \{1, \dots, m\}$ is the composition of a d -dimensional MIMO filter with a
 176 decimator. It follows from Noble identities [SM6] that \mathcal{W} reduces to cascading a $\zeta_m \times \zeta_0$ MIMO
 177 filter with frequency response $\widehat{\mathbf{W}}$ with a decimation of each output by a factor σ_m . More
 178 precisely, if $\mathbf{x} \in \mathcal{H}_0$ is the input of this linear system and \mathbf{y} its output, we have in the frequency
 179 domain:

$$180 \quad (\forall \boldsymbol{\nu} \in [0, 1]^d) \quad \widehat{\mathbf{y}}(\boldsymbol{\nu}) = \frac{1}{\sigma_m^d} \sum_{\mathbf{j} \in \mathbb{S}(\sigma_m)} \widehat{\mathbf{W}} \left(\frac{\boldsymbol{\nu} + \mathbf{j}}{\sigma_m} \right) \widehat{\mathbf{x}} \left(\frac{\boldsymbol{\nu} + \mathbf{j}}{\sigma_m} \right)$$

$$181 \quad (\text{SM4.5}) \quad = \frac{1}{\sigma_m^d} \widetilde{\mathbf{W}} \left(\frac{\boldsymbol{\nu}}{\sigma_m} \right) \widetilde{\mathbf{x}} \left(\frac{\boldsymbol{\nu}}{\sigma_m} \right),$$
 182

183 where $\widetilde{\mathbf{x}} \left(\frac{\boldsymbol{\nu}}{\sigma_m} \right)$ is a vector of dimension $\mathbb{C}^{\sigma_m^d \zeta_0}$ where the vectors $(\widehat{\mathbf{x}}((\boldsymbol{\nu} + \mathbf{j})/\sigma_m))_{\mathbf{j} \in \mathbb{S}(\sigma_m)}$
 184 are stacked columnwise and $\widetilde{\mathbf{W}} \left(\frac{\boldsymbol{\nu}}{\sigma_m} \right)$ is a $c_m \times \sigma_m^d \zeta_0$ matrix where the matrices $(\widehat{\mathbf{W}}((\boldsymbol{\nu} +$
 185 $\mathbf{j})/\sigma_m))_{\mathbf{j} \in \mathbb{S}(\sigma_m)}$ are stacked rowwise. For example, when $d = 2$, we have, for every $\boldsymbol{\nu} =$

² $(\cdot)^H$ denotes the Hermitian transpose operation.

186 $(\nu_1, \nu_2) \in [0, 1]^2$,

$$187 \quad (\text{SM4.6}) \quad \tilde{\mathbf{x}}(\boldsymbol{\nu}) = \begin{bmatrix} \check{\mathbf{x}}(\nu_1, \nu_2) \\ \check{\mathbf{x}}\left(\nu_1, \nu_2 + \frac{1}{\sigma_m}\right) \\ \vdots \\ \check{\mathbf{x}}\left(\nu_1, \nu_2 + \frac{\sigma_m - 1}{\sigma_m}\right) \end{bmatrix} \in \mathbb{C}^{\sigma_m^2 \zeta_0}$$

$$188 \quad (\text{SM4.7}) \quad \check{\mathbf{x}}(\boldsymbol{\nu}) = \begin{bmatrix} \hat{\mathbf{x}}(\nu_1, \nu_2) \\ \hat{\mathbf{x}}\left(\nu_1 + \frac{1}{\sigma_m}, \nu_2\right) \\ \vdots \\ \hat{\mathbf{x}}\left(\nu_1 + \frac{\sigma_m - 1}{\sigma_m}, \nu_2\right) \end{bmatrix} \in \mathbb{C}^{\sigma_m \zeta_0}$$

$$189 \quad (\text{SM4.8}) \quad \widetilde{\mathbf{W}}(\boldsymbol{\nu}) = \left[\check{\mathbf{W}}(\nu_1, \nu_2) \quad \check{\mathbf{W}}\left(\nu_1, \nu_2 + \frac{1}{\sigma_m}\right) \quad \dots \quad \check{\mathbf{W}}\left(\nu_1, \nu_2 + \frac{\sigma_m - 1}{\sigma_m}\right) \right] \in \mathbb{C}^{\zeta_m \times \sigma_m^2 \zeta_0}$$

$$190 \quad (\text{SM4.9}) \quad \check{\mathbf{W}}(\boldsymbol{\nu}) = \left[\widehat{\mathbf{W}}(\nu_1, \nu_2) \quad \widehat{\mathbf{W}}\left(\nu_1 + \frac{1}{\sigma_m}, \nu_2\right) \quad \dots \quad \widehat{\mathbf{W}}\left(\nu_1 + \frac{\sigma_m - 1}{\sigma_m}, \nu_2\right) \right] \in \mathbb{C}^{\zeta_m \times \sigma_m \zeta_0}.$$

192 By using now Parseval's formula,

$$\begin{aligned} 193 \quad \|\mathbf{y}\|^2 &= \int_{[0,1]^d} \|\widehat{\mathbf{y}}(\boldsymbol{\nu})\|^2 d\boldsymbol{\nu} \\ 194 &= \frac{1}{\sigma_m^{2d}} \int_{[0,1]^d} \left\| \widetilde{\mathbf{W}}\left(\frac{\boldsymbol{\nu}}{\sigma_m}\right) \tilde{\mathbf{x}}\left(\frac{\boldsymbol{\nu}}{\sigma_m}\right) \right\|^2 d\boldsymbol{\nu} \\ 195 &\leq \frac{1}{\sigma_m^d} \int_{[0,1/\sigma_m]^d} \left\| \widetilde{\mathbf{W}}(\boldsymbol{\nu}) \right\|^2 \|\tilde{\mathbf{x}}(\boldsymbol{\nu})\|^2 d\nu_1 d\nu_2 \\ 196 &\leq \frac{1}{\sigma_m^d} \sup_{\boldsymbol{\nu} \in [0,1/\sigma_m]^d} \left\| \widetilde{\mathbf{W}}(\boldsymbol{\nu}) \right\|^2 \int_{[0,1/\sigma_m]^d} \|\tilde{\mathbf{x}}(\boldsymbol{\nu})\|^2 d\boldsymbol{\nu} \\ 197 &= \frac{1}{\sigma_m^2} \sup_{\boldsymbol{\nu} \in [0,1/\sigma_m]^d} \left\| \widetilde{\mathbf{W}}(\boldsymbol{\nu}) \right\|^2 \sum_{\mathbf{j} \in \mathbb{S}(\sigma_m)} \int_{[0,1/\sigma_m]^d} \left\| \hat{\mathbf{x}}\left(\boldsymbol{\nu} + \frac{\mathbf{j}}{\sigma_m}\right) \right\|^2 d\boldsymbol{\nu} \\ 198 &= \frac{1}{\sigma_m^d} \sup_{\boldsymbol{\nu} \in [0,1/\sigma_m]^d} \left\| \widetilde{\mathbf{W}}(\boldsymbol{\nu}) \right\|^2 \int_{[0,1]^d} \|\hat{\mathbf{x}}(\boldsymbol{\nu})\|^2 d\boldsymbol{\nu} \\ 199 \quad (\text{SM4.10}) &= \frac{1}{\sigma_m^d} \sup_{\boldsymbol{\nu} \in [0,1/\sigma_m]^d} \left\| \widetilde{\mathbf{W}}(\boldsymbol{\nu}) \right\|^2 \|\mathbf{x}\|^2. \\ 200 \end{aligned}$$

201 This shows that

$$202 \quad (\text{SM4.11}) \quad \theta_m^2 \leq \frac{1}{\sigma_m^d} \sup_{\boldsymbol{\nu} \in [0,1/\sigma_m]^d} \left\| \widetilde{\mathbf{W}}(\boldsymbol{\nu}) \right\|^2.$$

203 On the other hand since $\widehat{\mathbf{W}}$ is continuous, $\widetilde{\mathbf{W}}$ is also continuous, and there exists $\hat{\boldsymbol{\nu}} \in [0, 1/\sigma_m]^d$
204 such that

$$205 \quad (\text{SM4.12}) \quad \sup_{\boldsymbol{\nu} \in [0,1/\sigma_m]^d} \left\| \widetilde{\mathbf{W}}(\boldsymbol{\nu}) \right\| = \left\| \widetilde{\mathbf{W}}(\hat{\boldsymbol{\nu}}) \right\|.$$

206 Let us now choose, for every $\boldsymbol{\nu} \in [0, 1/\sigma_m]^d$, $\tilde{\mathbf{x}}(\boldsymbol{\nu}) = \alpha_\epsilon(\boldsymbol{\nu})\mathbf{u}(\boldsymbol{\nu})$ where $\mathbf{u}(\boldsymbol{\nu})$ is a unit norm
 207 eigenvector associated with the maximum eigenvalue of $\widetilde{\mathbf{W}}(\boldsymbol{\nu})^H \widetilde{\mathbf{W}}(\boldsymbol{\nu})$, $\epsilon \in]0, +\infty[$, and

$$208 \quad (\text{SM4.13}) \quad \alpha_\epsilon(\boldsymbol{\nu}) = \begin{cases} \frac{1}{\epsilon^{d/2}} & \text{if } (\exists \mathbf{j} \in \{-1, 0, 1\}^d) \|\boldsymbol{\nu} + \frac{\mathbf{j}}{\sigma_m} - \widehat{\boldsymbol{\nu}}\|_\infty \leq \frac{\epsilon}{2} \\ 0 & \text{otherwise.} \end{cases}$$

209 Then we see that when $\epsilon \rightarrow 0$, the upper bound in (SM4.10) is reached. We conclude that

$$210 \quad (\text{SM4.14}) \quad \theta_m = \frac{1}{\sigma_m^{d/2}} \sup_{\boldsymbol{\nu} \in [0, 1/\sigma_m]^d} \|\widetilde{\mathbf{W}}(\boldsymbol{\nu})\|.$$

211 In addition, by using the relation between $\widetilde{\mathbf{W}}$ and $\widehat{\mathbf{W}}$ (i.e., (SM4.8) and (SM4.9) in the 2D
 212 case),

$$213 \quad \|\widetilde{\mathbf{W}}(\boldsymbol{\nu})\|^2 = \|\widetilde{\mathbf{W}}(\boldsymbol{\nu})\widetilde{\mathbf{W}}(\boldsymbol{\nu})^H\| \\ 214 \quad (\text{SM4.15}) \quad = \left\| \sum_{\mathbf{j} \in \mathbb{S}(\sigma_m)} \widehat{\mathbf{W}}\left(\boldsymbol{\nu} + \frac{\mathbf{j}}{\sigma_m}\right) \widehat{\mathbf{W}}\left(\boldsymbol{\nu} + \frac{\mathbf{j}}{\sigma_m}\right)^H \right\|. \\ 215$$

216 Gathering the last two equalities yields (SM4.2). ■

217 When there is no decimation, i.e. the strides $(s_i)_{1 \leq i \leq m}$ are all equal to 1, (SM4.2) reduces
 218 to

$$219 \quad (\text{SM4.16}) \quad \theta_m = \sup_{\boldsymbol{\nu} \in [0, 1]^d} \|\widehat{\mathbf{W}}_m(\boldsymbol{\nu}) \cdots \widehat{\mathbf{W}}_2(\boldsymbol{\nu}) \widehat{\mathbf{W}}_1(\boldsymbol{\nu})\|.$$

220 We recall that the following upper bound holds [SM49]:

$$221 \quad (\text{SM4.17}) \quad \theta_m \leq \bar{\theta}_m = \prod_{i=1}^m \|\mathcal{W}_i\|.$$

222 Applying our result to the one-layer case shows that, for every $i \in \{1, \dots, m\}$,

$$223 \quad (\text{SM4.18}) \quad \|\mathcal{W}_i\| = \frac{1}{s_i^{d/2}} \sup_{\boldsymbol{\nu} \in [0, 1/s_i]^2} \left\| \sum_{\mathbf{j} \in \mathbb{S}(s_i)} \widehat{\mathbf{W}}_i\left(\boldsymbol{\nu} + \frac{\mathbf{j}}{s_i}\right) \widehat{\mathbf{W}}_i\left(\boldsymbol{\nu} + \frac{\mathbf{j}}{s_i}\right)^H \right\|^{1/2}.$$

224 Note that the resulting upper bound in (SM4.17) gives a loose estimate of the Lipschitz
 225 constant, which has however the merit to be valid for convolutional networks having kernels
 226 with an arbitrary sign.

227 **SM5. Numerical evaluation of the Lipschitz constant of nonnegative convolutional networks.**
 228 We compare the tight bound θ_m in Theorem 4.1 with the separable one $\bar{\theta}_m$ given by (4.12) for
 229 a classic convolutional network using non-negative kernels. The results provided in Table SM1
 230 correspond to the convolutive part of LeNet-5 [SM4]. In our experiments, we initialized

231 the networks with randomly sampled weights drawn from a uniform distribution on $[0, 1]$.
 232 Table SM1 shows the relative difference

$$233 \quad \epsilon_r = \frac{\bar{\theta}_m - \theta_m}{\bar{\theta}_m},$$

234 for 10 distinct noise realizations. We thus observe that the difference between the two bounds
 235 is small. Similar observations can be made on various convolutive architectures. In contrast,
 236 for fully connected networks, a separable bound is usually overpessimistic.

LeNet-5										
	#1	#2	#3	#4	#5	#6	#7	#8	#9	#10
θ_m	30302.73	27734.91	30298.73	29374.35	30180.16	28632.60	30615.02	30395.67	34828.90	30097.62
$\bar{\theta}_m$	30696.07	28114.29	30860.56	29821.62	30670.05	29298.64	31152.06	30866.87	35220.36	30367.71
ϵ_r [%]	1.28	1.35	1.82	1.50	1.60	2.27	1.72	1.53	1.11	0.89

Table SM1: Lipschitz bounds obtained for 10 independent realizations of random positive initialization for LeNet-5.

237 **SM6. Lipschitz constant of average pooling.** We consider the case when the i -th layer is
 238 an *average pooling* where the average is computed on patches of length L_i in each dimension
 239 and with stride s_i . For simplicity, we suppose that L_i is a multiple of s_i . The number of input
 240 and output channels is then equal, i.e. $\zeta_i = \zeta_{i-1}$. The average is calculated on each channel
 241 independently, this operation is a special case of a nonnegative convolutional layer where, for
 242 every $\mathbf{n} \in \mathbb{Z}^d$, $\mathbf{W}_i(\mathbf{n})$ is a diagonal matrix. The diagonal elements of this matrix are

$$243 \quad (\text{SM6.1}) \quad (\forall p \in \{1, \dots, \zeta_i\})(\forall \mathbf{n} \in \mathbb{Z}^d) \quad w_{i,p,p}(\mathbf{n}) = \begin{cases} \frac{1}{L_i^d} & \text{if } \mathbf{n} \in [0, L_i - 1]^d \\ 0 & \text{otherwise.} \end{cases}$$

244 We deduce that, for every $\mathbf{j} \in \mathbb{S}(s_i)$, the matrix $\bar{\mathbf{W}}_i^{(\mathbf{j})}$ is also a diagonal matrix. More precisely,
 245 the sum in (4.13) can be restricted to values of $\mathbf{n} \in \{0, \dots, L_i/s_i - 1\}^d$ and $\bar{\mathbf{W}}_i^{(\mathbf{j})} = \frac{1}{s_i^d} \mathbf{Id}$. We
 246 deduce that the Lipschitz constant of the average pooling layer is

$$247 \quad (\text{SM6.2}) \quad \|\mathcal{W}_i\| = \left\| \sum_{\mathbf{j} \in \mathbb{S}(s_i)} \bar{\mathbf{W}}_i^{(\mathbf{j})} (\bar{\mathbf{W}}_i^{(\mathbf{j})})^\top \right\|^{1/2} = \frac{1}{s_i^{d/2}}.$$

248 We see that this constant is independent of the patch size and is a decreasing function of the
 249 stride.

250 **SM7. Expressivity of ABBA networks – simulations.** For this experiment, we randomly sam-
 251 pled points from four distinct 2D Gaussian distributions, with different means and covariance
 252 matrices, totaling 125 2-dimensional points per class. Figure SM1 shows a comparison between
 253 decision boundaries resulting from training two models: a standard one trained conventionally
 254 and its non-negative ABBA equivalent. The two models reach a similar solution, showing that
 255 the theoretical properties proved in this paper are also observed in practical simulations.

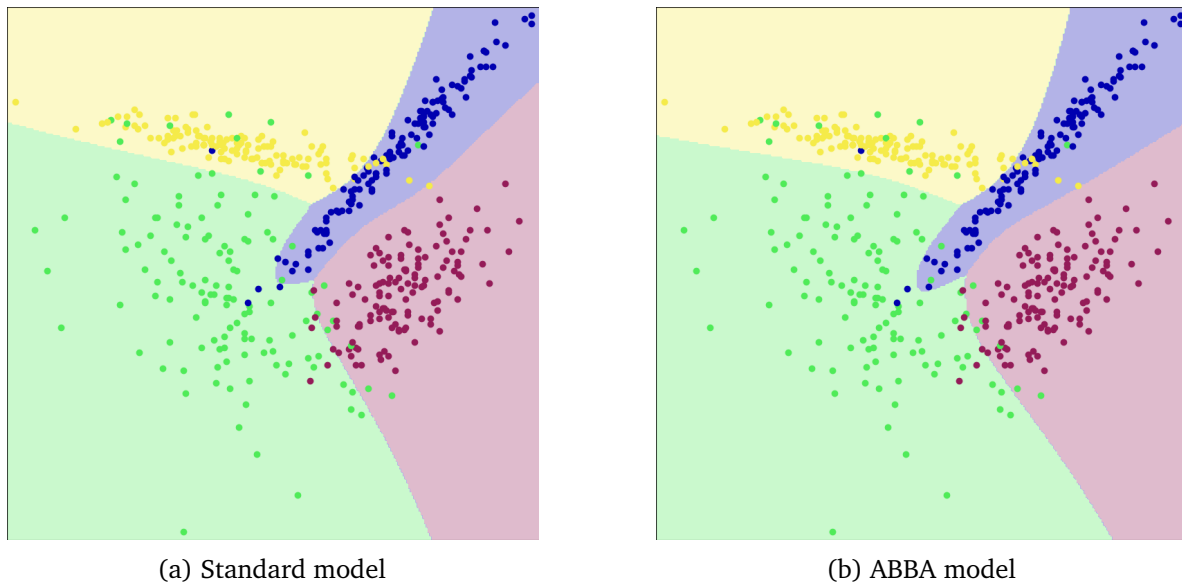


Figure SM1: Decision space comparison between fitting an ABBA network and a standard arbitrary-signed one.

256 **SM8. Constrained training of signed convolutional layers.** The first and the last layers of an
 257 ABBA convolutional network have signed kernels. The norm of these layers is computed by
 258 using (SM4.18) and constrained to be less than $\bar{\theta}_{m,i,t}$ with $i \in \{0, m + 1\}$. Note that (SM4.18)
 259 makes use of the frequency response $\widehat{\mathbf{W}}_i$ of filter \mathcal{W}_i . A discrete Fourier transform (DFT) is
 260 actually implemented (using 128×128 discrete frequencies). In the discrete frequency domain,
 261 the upper bound constraint is thus decomposed into 128^2 matrix norm bounds obtained by
 262 summing over s_i^2 frequencies. The projection onto each of these elementary constraint sets is
 263 computed by truncating a singular value decomposition. An additional constraint, however, is
 264 to be addressed, which is related to the fact that the kernels are of finite size. This implicitly
 265 defines a linear constraint. Projecting onto the associated vector space is simply obtained by
 266 truncating the kernel (after inverse DFT) to the desired size. The set $\mathcal{S}_{i,t}$ is thus defined as the
 267 intersection of the former matrix norm constraint set and the latter vector space. Projecting
 268 onto this intersection can be achieved by an iterative convex optimization approach. In our
 269 case, we use a Douglas-Rachford algorithm [SM2].

270 **SM9. ABBA architectures.** Table SM3 details the ABBA Dense and ABBA Conv architectures
 271 used for MNIST and FMNIST datasets, while Table SM2 shows our choices for RPS and CelebA
 272 datasets. As the ABBA layers have a specific form, their output size will be twice the number
 273 of filters. The used activation operator is the Capped Leaky ReLu (CLR) function defined in
 274 (SM1.3) for all Dense layers. For convolutional operators we employed a 3×3 kernel, using
 275 the same activation.

276 We used the official train-test split provided by the Tensorflow framework for both MNIST

Layer type	RPS	stride	CelebA	stride
Input	$150 \times 150 \times 3$		$128 \times 128 \times 3$	
Conv2D	$150 \times 150 \times 8$	1	$128 \times 128 \times 8$	1
ABBA Conv2D + CLR	–	–	$64 \times 64 \times 8(\times 2)$	2
ABBA Conv2D + CLR	$75 \times 75 \times 32(\times 2)$	2	$32 \times 32 \times 16(\times 2)$	2
ABBA Conv2D + CLR	$37 \times 37 \times 64(\times 2)$	2	$16 \times 16 \times 32(\times 2)$	2
ABBA Conv2D + CLR	$18 \times 18 \times 128(\times 2)$	2	$8 \times 8 \times 64(\times 2)$	2
Conv2D	$18 \times 18 \times 128$	1	$8 \times 8 \times 64$	1
Global Max-Pooling2D	$128(\times 2)$		$64(\times 2)$	
ABBA Dense + CLR	$128(\times 2)$		–	
ABBA Dense + CLR	$64(\times 2)$		–	
ABBA Dense + CLR	$32(\times 2)$		–	
Dense	3		2	

Table SM2: ABBA Conv architectures details for RPS and CelebA datasets.

Layer type	MNIST/FMNIST	Layer type	MNIST	FMNIST
Input	$28 \times 28 \times 1$	Input	784	784
Conv2D	$28 \times 28 \times 32$	Dense	256	256
ABBA Conv2D + CLR	$28 \times 28 \times 16$	ABBA Dense + CLR	128	128
ABBA Conv2D + CLR	$28 \times 28 \times 16$	ABBA Dense + CLR	64	64
Conv2D	$28 \times 28 \times 1$	ABBA Dense + CLR	–	32
Dense	256	Dense	10	10
ABBA Dense + CLR	128			
ABBA Dense + CLR	64			
Dense	10			

Table SM3: ABBA Dense and ABBA Conv architecture details for MNIST and FMNIST datasets. For convolutional layers the stride is set to 1.

277 and FMNIST datasets and did not employ any augmentation strategy during training. For RPS
 278 and CelebA models, we resized the input images to 150×150 , resp. 128×128 , before feeding
 279 them to the network. In the case of CelebA dataset, we opted for a binary classification task on
 280 the *bald* feature. We extracted all the images containing the *bald* attribute, and we randomly
 281 select the same number of examples from the *non-bald* class, in order to avoid class imbalance.
 282 Additional information regarding the optimization parameters used during training is provided
 283 in Table SM4.

284 **SM10. Adversarial examples.** For all datasets, adversarial examples created by using DDN
 285 attack are displayed in Figures SM2, SM3, SM4, and SM5. We generated adversarial samples
 286 using untargeted DDN attacks, with a budget of 300 iterations and initial parameters as
 287 proposed by the authors. We did not limit the maximum perturbation ϵ , in order to find the
 288 minimum one allowing us to fool the model. It can be easily seen that for Deelip and ABBA



Figure SM2: Adversarial examples with DDN attack for Conv-Dense models, on MNIST dataset. l_2 perturbation magnitude is given in the top-left corner.

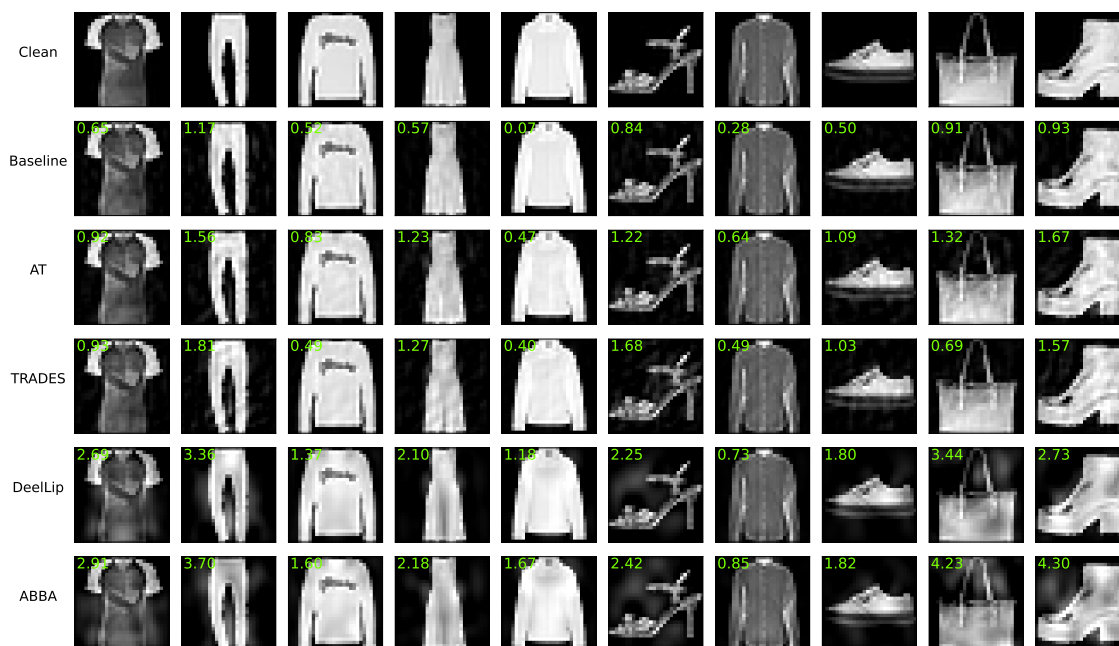


Figure SM3: Adversarial examples with DDN attack for Conv-Dense models, on FMNIST dataset. l_2 perturbation magnitude is given in the top-left corner.

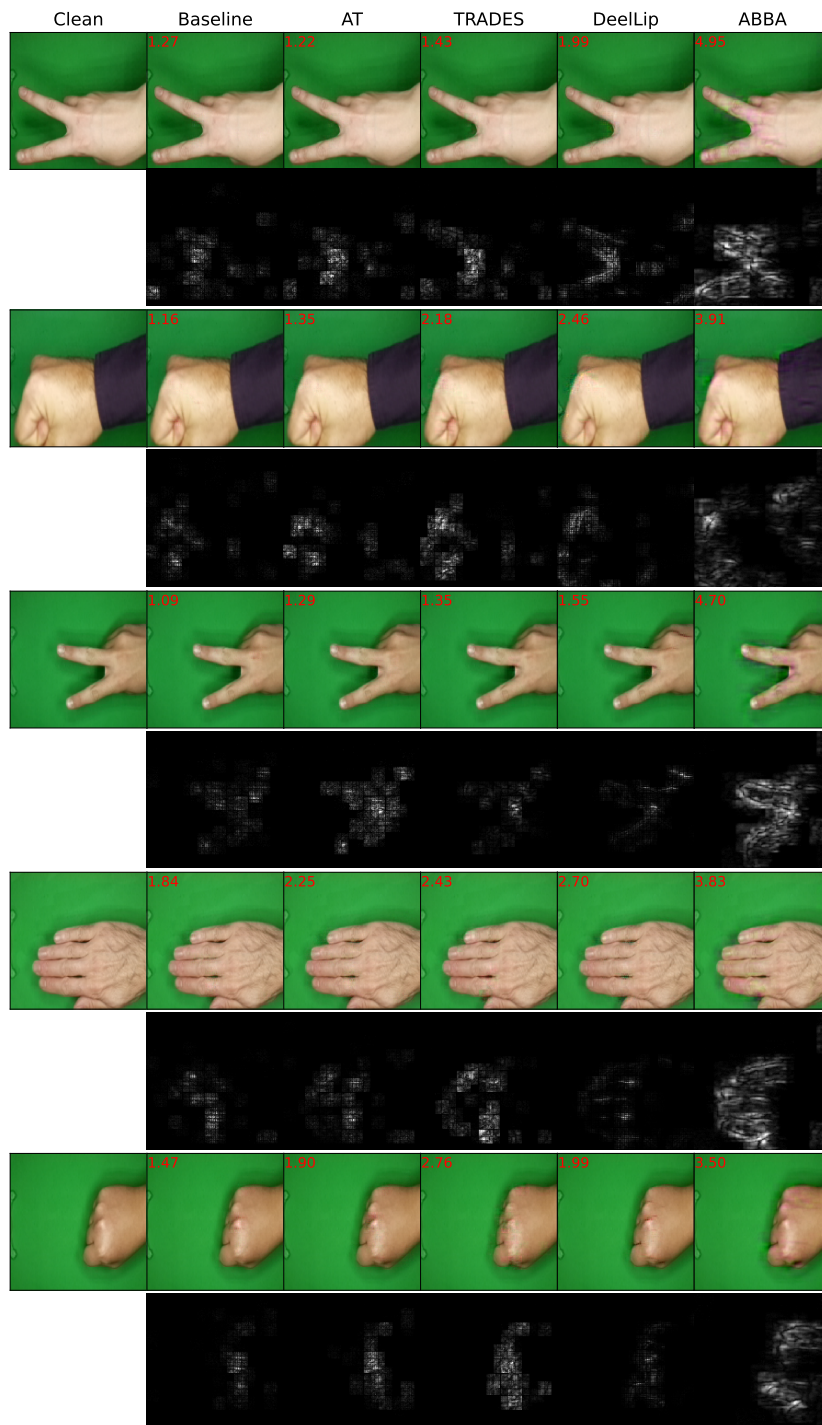


Figure SM4: Adversarial examples generated with DDN, on RPS dataset. For each example: first row – adversarial images; second row – pixel differences between adversarial and clean sample.



Figure SM5: Adversarial examples with DeepFool attack for CelebA. l_2 perturbation magnitude is given in the top-left corner.

Dataset	Optimizer	No. Epochs	Learning rate	Batch size
MNIST	projected ADAM	150	10^{-3}	1024
FMNIST	projected ADAM	200	10^{-3}	1024
RPS	projected ADAM	250	10^{-4}	64
CelebA	projected ADAM	100	10^{-4}	128

Table SM4: Training hyperparameters.

289 networks the required perturbations for misclassification are higher. In particular, we observe
 290 that the perturbations needed to fool ABBA networks lead to severe artifacts in the images.

291 **SM11. Training time.** We first compared the average time/epoch for training a standard
 292 network and its ABBA equivalent. Table SM5 reports the average seconds per epoch for both
 293 cases, for different feed-forward architectures. On average, training an ABBA neural network
 294 for 200 epochs on MNIST introduces less than 10% additional training time.

295 The projection is a costly step, and it is the main contributor to the training overhead. A
 296 comparison of the training time (per-batch) with (*green line*) and without (*dotted green line*)
 297 projection is featured in Figures SM6a and SM6b for architectures with an increasing number
 298 of fully connected and convolutional layers, respectively. The average deviation from the
 299 imposed global bound (*red*), which was set to 1 in all cases, is also reported. This shows that
 300 we are able to maintain the imposed bounds, given the same number of iterations, irrespective
 301 of the network depth.

Model		Architecture					
		2C2F	2C3F	3C2F	4C2F	4C3F	5C1F
Standard	Acc [%]	95.28	95.80	99.18	99.30	99.26	99.10
	Sec./Epoch	4.25	4.29	4.31	4.28	4.37	4.31
	Size (MB)	0.09	0.14	0.39	0.53	0.59	0.97
ABBA	Acc [%]	95.54	95.34	98.62	99.12	99.14	98.72
	Sec./Epoch	4.52	4.61	4.67	4.68	4.79	4.74
	Size (MB)	0.18	0.28	0.78	1.06	1.18	1.94

Table SM5: Comparison of per-epoch training times for various Standard and ABBA architectures, on MNIST. *XCYP* corresponds to an architecture with *X* Convolution layers, followed by *Y* fully-connected layers.

302

REFERENCES

303 [1] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert spaces*, 2nd
 304 ed., corrected printing. New York: Springer, (2019).

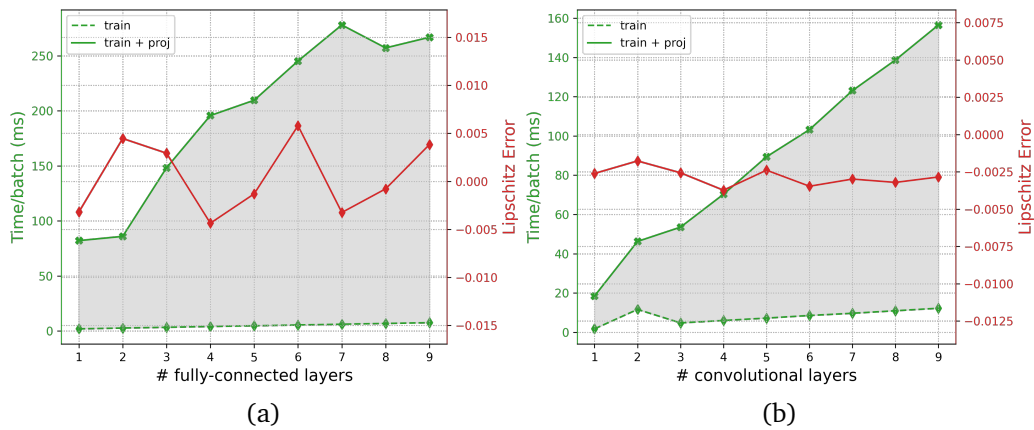


Figure SM6: Computation time for the projection step for a variable-length sequence of ABBA **SM6a** fully-connected and **SM6b** convolutional layers. All projections were computed with a number of 10 iterations, and the results were averaged over 50 independent simulations.

305 [2] P. L. COMBETTES AND J.-C. PESQUET, *Proximal splitting methods in signal processing*, Fixed-Point Algorithms
 306 for Inverse Probl. Sci. Eng., (2011), pp. 185–212.
 307 [3] ———, *Lipschitz certificates for layered network structures driven by averaged activation operators*, J. Math.
 308 Data Sci., 2 (2020), pp. 529–557.
 309 [4] Y. LECUN, L. BOTTOU, Y. BENGIO, AND P. HAFFNER, *Gradient-based learning applied to document recognition*,
 310 Proc. IEEE, 86 (1998), pp. 2278–2324.
 311 [5] J. R. SILVESTER, *Determinants of block matrices*, The Math. Gazette, 84 (2000), pp. 460–467.
 312 [6] P. P. VAIDYANATHAN, *Multirate Systems and Filter Banks*, Prentice Hall (NJ), (1993).