

## Basics of Neural Networks

DEEP LEARNING IS EVERYWHERE

VISION

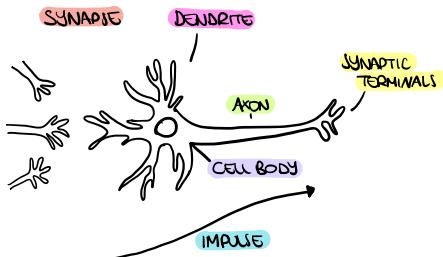
LANGUAGE

SPEECH

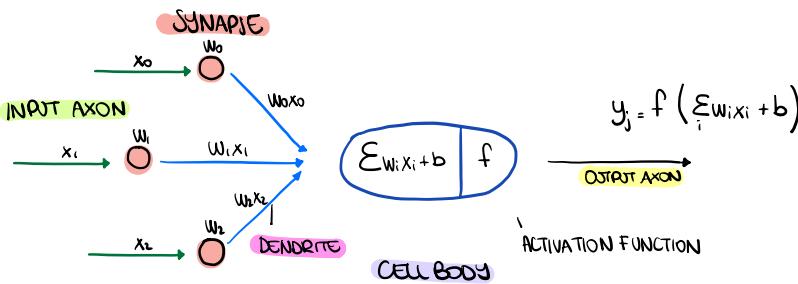
THE PERFORMANCE FOR DEEP LEARNING IS IMPROVING VERY RAPIDLY.

## Terminology of Neural Networks

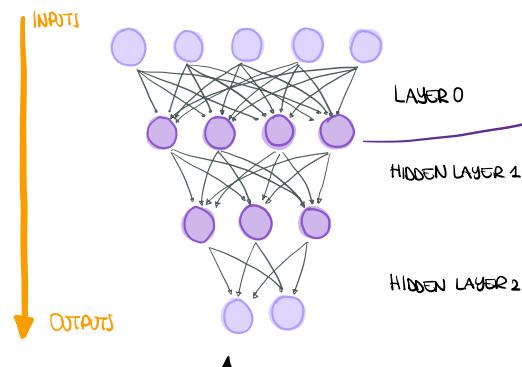
### Neuron and Synapse



A synapse has many input axons and then we multiply the weights of the input axons and combine them together.



### Example of a 3-layer neural network



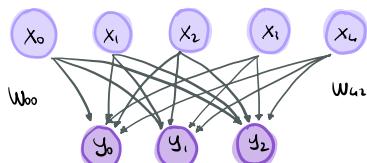
#### Note on Terminology:

SYNAPSES = WEIGHTS = PARAMETERS

NEURONS = FEATURES = ACTIVATIONS

THE DIMENSIONALITY OF THESE HIDDEN LAYERS DETERMINES THE WIDTH OF THE MODEL

THIS NETWORK HAS FULLY-CONNECTED LAYERS (LINEAR LAYERS) → THE OUTPUT NEURON IS CONNECTED TO ALL INPUT NEURONS.



$$y_i = \sum_j w_{ij} x_j + b_i$$

#### Shape of Tensors:

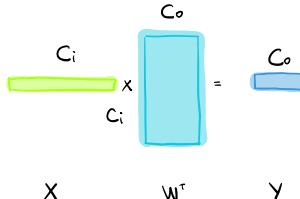
INPUT FEATURES  $X$ :  $(1, c_i)$

OUTPUT FEATURES  $Y$ :  $(1, c_o)$

WEIGHTS  $W$ :  $(c_o, c_i)$

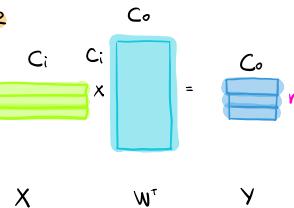
BIAS  $b$ :  $(c_o,)$

WHERE  $c_i$  ARE THE INPUT CHANNELS AND  $c_o$  THE OUTPUT CHANNELS.



HOW TO CALCULATE THE DIMENSION OF A MULTILAYER PERCEPTRON:

IF WE HAVE A BATCH SIZE OF  $n$  →



#### Shape of Tensors:

INPUT FEATURES  $X$ :  $(n, c_i)$

OUTPUT FEATURES  $Y$ :  $(n, c_o)$

WEIGHTS  $W$ :  $(c_o, c_i)$

BIAS  $b$ :  $(c_o,)$

LET'S NOW LOOK AT A SLIGHTLY MORE COMPLICATED SCENARIO → CONVOLUTION LAYER



THE OUTPUT NEURON IS CONNECTED TO INPUT NEURONS IN THE RECEPTIVE FIELD



RATHER THAN HAVING JUST ONE CHANNEL, WE HAVE MULTIPLE DIMENSIONS IN THIS SPATIAL DIMENSION.

MANY SIGNALS ARE SIMILAR TO THIS CASE E.G.: SPEECH SIGNALS

- FOR SPEECHrecognition WE HAVE A LONG TEMPORAL SEQUENCE

#### Shape of Tensors:

1D CONV  
INPUT FEATURES  $X$ :  $(n, c_i, w_i)$

OUTPUT FEATURES  $Y$ :  $(n, c_o, w_o)$

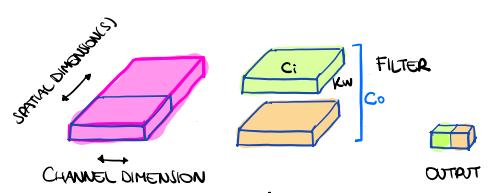
WEIGHTS  $W$ :  $(c_o, c_i, k_w)$

KERNEL WIDTH

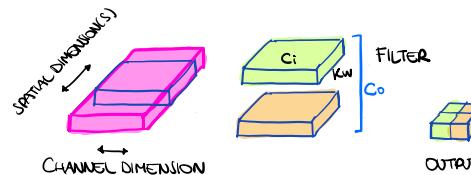
BIAIS  $b$ :  $(c_o,)$

THE ORANGE FILTER PRODUCES THE ORANGE OUTPUT AND THE GREEN FILTER PRODUCES THE GREEN OUTPUT.

THE ONE-DIMENSIONAL CONVOLUTION IS COMPUTED BY PERFORMING THE MULTIPLICATION OF EACH ELEMENT AND SUM THEM UP TOGETHER TO PRODUCE ONE OUTPUT ELEMENT. THERE ARE MULTIPLE FILTERS WHICH CAN PRODUCE ANOTHER CHANNEL IN THE SAME DIMENSION.



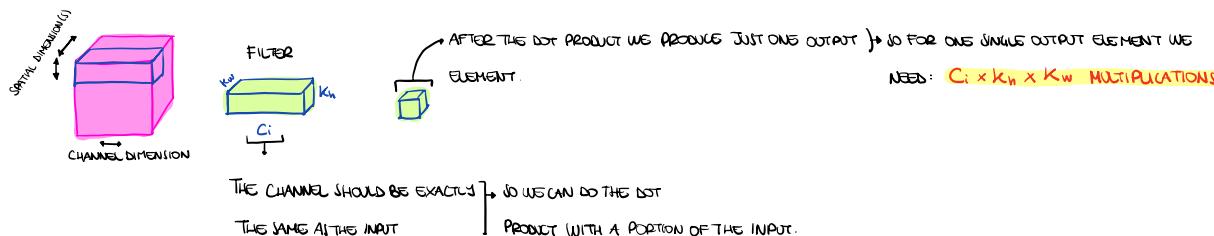
IN ORDER TO PRODUCE THE NEXT SPATIAL DIMENSION WE CAN SHIFT THE INPUT BY ONE → SO WE CAN PRODUCE THE NEXT SPATIAL OUTPUT.



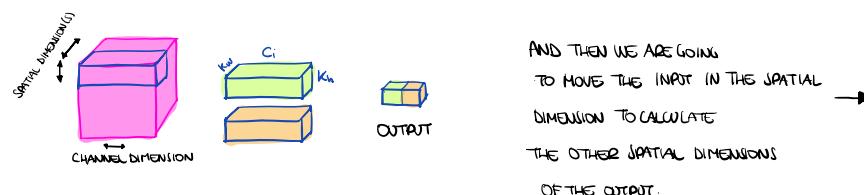
LET'S MOVE TO A MORE COMPLICATED SCENARIO → 2D CONVOLUTION → WIDELY USED IN PROCESSING



IN ORDER TO PRODUCE THE OUTPUT WE USE AGAIN A FILTER → THE DIMENSION OF THE FILTER IS DENOTED BY THE K, WHICH IS INVOKING "KERNEL".

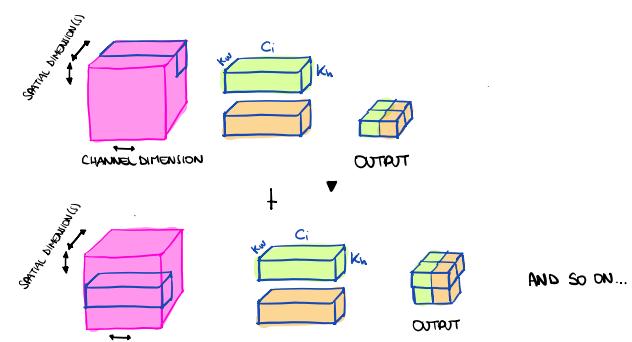


IF WE WANT TO HAVE MULTIPLE OUTPUT CHANNELS WE ARE GOING TO HAVE MULTIPLE FILTERS:



HOW DO WE CALCULATE THE OUTPUT DIMENSION?

↪ OUTPUT FEATURES Y →  $(n, C_o, h_o, w_o)$   
WHERE  $h_o = h_i - K_h + 1$



**PADDING** → PADDING CAN BE USED TO KEEP THE OUTPUT FEATURE MAP SIZE THE SAME AS THE INPUT FEATURE MAP SIZE.

THERE ARE SEVERAL APPROACHES TO DO THE PADDING, THE SIMPLEST ONE IS TO DO ZERO PADDING.

↪ WITH PADDING THE OUTPUT SIZE IS CALCULATED USING:  $h_o = h_i + 2p - K_h + 1$

↪ WHERE p IS THE PADDING SIZE

OTHER PADDING:

- REFLECTION PADDING → IT "REFLECTS" THE ROW INTO THE PADDING

- REPETITION PADDING → REPETITION IS DONE BY COPYING THE NEAREST BORDER VALUES

0	0	0	0	0	0	0
0	0	0	0	0	0	0
0	0	1	2	3	0	0
0	0	4	5	6	0	0
0	0	7	8	9	0	0
0	0	0	0	0	0	0
0	0	0	0	0	0	0

IT PADS THE INPUT BOUNDARIES WITH ZERO (IT'S THE DEFAULT IN PYTORCH)

**CONVOLUTION LAYER → RECEPTIVE FIELD** → WE WANT TO HAVE A GLOBAL UNDERSTANDING OF THE WHOLE IMAGE TO UNDERSTAND WHAT ARE THE HIGH LEVEL FEATURES.

↪ IN CONVOLUTION EACH OUTPUT ELEMENT DEPENDS ON  $K_h \times K_w$  RECEPTIVE FIELD IN THE INPUT.

• EACH SUBSIDIARY CONVOLUTION ADDS  $K-1$  TO THE RECEPTIVE FIELD SIZE

• WITH L LAYERS, THE RECEPTIVE FIELD SIZE IS  $L \cdot (K-1) + 1$

↪ THE PROBLEM IS THAT IN ORDER TO HAVE A LARGE RECEPTIVE FIELD WE NEED MANY LAYERS FOR EACH OUTPUT TO "SEE" THE WHOLE IMAGE.

SOLUTION: DOWNSAMPLE INSIDE THE NEURAL NETWORK

↪ A METHOD TO REDUCE THE AMOUNT OF WEIGHTS

**GROUPED CONVOLUTION LAYER** → A GROUP OF NARROWER CONVOLUTIONS

SHAPE OF TENSORS

INPUT FEATURES X:  $(n, C_i, h_i, w_i)$

OUTPUT FEATURES Y:  $(n, C_o, h_o, w_o)$

WEIGHTS W:  $(g \cdot C_o / g, C_i / g, K_h, K_w)$

BIAIS b:  $(C_o)$

GROUDED

PREVIOUSLY EACH OUTPUT CHANNEL WAS CALCULATED FOR EACH INPUT CHANNEL,

IN A GROUPED CONVOLUTION WE HAVE DIFFERENT GROUPS, EACH GROUP ONLY DEPENDS ON A SPECIFIC GROUP OF INPUT

RATHER THAN THE ENTIRE INPUT CHANNELS.

THE STRIDE CONVOLUTION ←  
LAYER INCREASES THE  
RECEPTIVE FIELD.

$$h_o = h_i + 2p - K_h + 1$$

S WHERE p IS PADDING AND S IS STRIDE.

STRIDED  
CONDUCTION  
LAYER

**DEPTHWISE CONVOLUTION LAYER** → THAT IS AN EXTREME CASE OF GROUPED CONVOLUTION WHERE THE NUMBER OF GROUPS IS EQUAL TO THE NUMBER OF CHANNELS.

INDEPENDENT FILTER FOR EACH CHANNEL →  $g = c_i = c_o$  IN GROUPED CONVOLUTION.

### SHAPE OF TENSORS

INPUT FEATURES X:  $(n, c_i, h_i, w_i)$

$$p_o = \frac{h_i + 2p - k_h + 1}{s}$$

OUTPUT FEATURES Y:  $(n, c_o, p_o, w_o)$

WEIGHTS W:  $(c, k_h, k_w)$

BIAIS b  $(c_o)$

MORE FANCY OPERATIONS:

DILATED CONVOLUTION   TRANSPOSED CONVOLUTION   ETC..

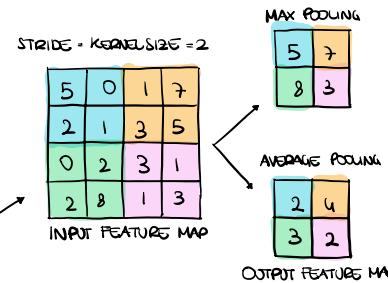
↳ USED FOR PIXELWISE PREDICTIONS

WHEN WE WANT TO DOWNSAMPLE THE FEATURE MAP TO A SMALLER SIZE

### POOLING LAYER

THE OUTPUT NEURON POOLS

USUALLY THE STRIDE IS THE SAME



THE FEATURES IN THE

AS THE KERNEL SIZE:  $s = k$

POOLING OPERATES OVER EACH CHANNEL INDEPENDENTLY

NO LEARNABLE PARAMETERS

RECEPTIVE FIELD

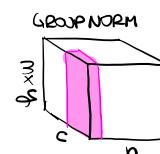
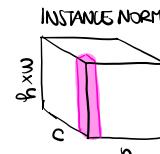
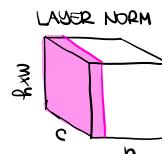
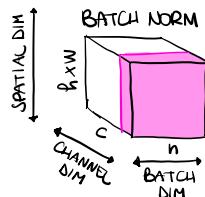
### NORMALIZATION LAYER

NORMALIZING THE FEATURES MAKES TRAINING FASTER

NORMALIZATION LAYER NORMALIZES THE FEATURES AS FOLLOWS →  $\hat{x}_i = \frac{1}{\sigma} (x_i - \mu)$  WHERE  $\mu_i = \frac{1}{m} \sum_{k \in S_i} x_k$  AND  $\sigma_i = \sqrt{\frac{1}{m} \sum_{k \in S_i} (x_k - \mu_i)^2 + \delta}$  ARE THE MEAN AND STANDARD DEVIATION OVER THE SET OF PIXELS  $S_i$

THEN LEARNS A PER-CHANNEL LINEAR TRANSFORM (TRAINABLE SCALE  $\gamma$  AND SHIFT  $\beta$ ) TO COMPENSATE FOR THE POSSIBLE LOSS OF REPRESENTATIONAL ABILITY →  $y = \gamma_i \hat{x}_i + \beta_i$

NOTE: DIFFERENT NORMALIZATIONS USE DIFFERENT DEFINITIONS OF THE SET  $S_i$  (COLORED IN PINK).



SOMETIMES THESE NORMALIZATION LAYERS CAN IMPROVE THE TRAINING CONVERGENCE BUT SOMETIMES THEY MAKE INFERENCE MORE CHALLENGING.

IN EFFICIENT NEURAL NETWORKS SOME OF THE NORMALIZATION LAYERS CAN BE FOLED INTO A CONVOLUTION LAYER TO SAVE THE COMPUTATION!

### ACTIVATION FUNCTION

ACTIVATION FUNCTIONS ARE TYPICALLY NON-LINEAR FUNCTIONS

**SIGMOID**

**ReLU**

**ReLU6**

**LEAKY RELU**

IT HELPS TO HAVE A BOOST IN THE ACCURACY BUT NOT VERY HARDWARE FRIENDLY FOR ITS IMPLEMENTATION  
SWISH  
 $y = \frac{x}{(1+e^{-x})}$   
NOT VERY QUANTIZATION FRIENDLY  
HARD SWISH  
 $y = x \cdot \frac{\text{ReLU6}(x+3)}{6}$

ONE OF THE MOST  
EFFICIENT FROM THE  
IMPLEMENTATION PERSPECTIVE

SINCE IT DOESN'T REQUIRE COMPLICATED ARITHMETIC

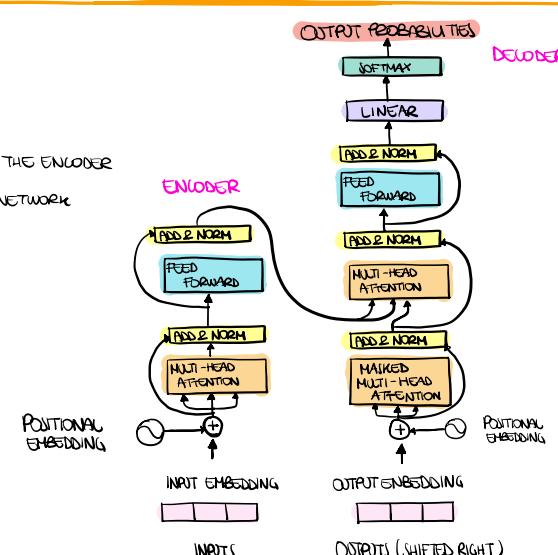
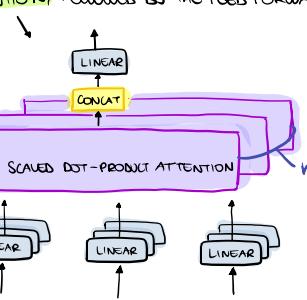
AND IT'S VERY QUANTIZATION FRIENDLY! → WE CAN THRESHOLD THE OUTPUT TO BE SOME NUMBER

**TRANSFORMER** THE TRANSFORMER HAS TWO PARTS: THE ENCODER AND THE DECODER. FOR THE ENCODER

WE HAVE TWO PARTS ONE IS THE MULTI-HEAD ATTENTION, FOLLOWED BY THE FEED FORWARD NETWORK

WE HAVE THREE INPUTS: THE VALUE, THE QUERY AND THE KEY. EACH OF THEM PASSES THROUGH A LINEAR TRANSFORMATION AND THEN WE DO THE ATTENTION MECHANISM → IT TRIES TO FIND THE RELATIONSHIP BETWEEN EACH TOKEN

$$\text{ATTENTION}(Q, K, V) = \text{SOFTMAX}\left(\frac{QK^T}{\sqrt{d_h}}\right)V$$



PROBLEM IF YOU HAVE A LONG SEQUENCE (WHERE  $n$  IS LARGE) IS A PROBLEM BECAUSE THE ATTENTION MAP IS  $n^2 \leftarrow$  THE COMPLEXITY GROWS QUADRATICALLY WITH THE SEQUENCE LENGTH

## POPULAR CNN ARCHITECTURES

### ALEXNET

• 5 CONVOLUTION LAYERS  
• 3 FULLY CONNECTED LAYERS  
AFTER EACH STAGE WITH A STRIDE OF 2 WE HAVE A RESOLUTION ROUGHLY HALF OF THE PREVIOUS STAGE

### VGG-16

IT HAS 16 LAYERS. EACH LAYER CONTAINS THE CONVOLUTION, BATCH NORMALIZATION AND THE RELU

### ResNet-50

IT INTRODUCED THE RESIDUAL BRANCH SO THAT WE CAN MAKE THE FLOW EASIER AND ENABLE A DEEPER NEURAL NET.

### MOBILE NETV2

USED FOR MOBILE APPLICATIONS  
↑ STEP-WISE CONVOLUTION

EFFICIENCY METRICS → HOW TO MEASURE THE EFFICIENCY OF DIFFERENT NEURAL NETS

FOR A NEURAL NET  
WE WANT IT TO BE SMALLER, FASTER AND GREENER → ENERGY CONSUMPTION  
STORAGE: HOW MANY STORAGE DO YOU NEED TO STORE YOUR NN  
LATENCY: HOW FAST CAN YOU PROCESS FROM THE INPUT TO THE OUTPUT

) ALL THESE CONCERN TWO FACTORS OF A NEURAL NETWORK: COMPUTATION AND MEMORY

NOTE:

MEMORY ≠ STORAGE. STORAGE IS A STATIC CONCEPT OF ONLY THE WEIGHTS AND CODE, MEMORY ALSO INCLUDES THE ACTIVATIONS.

THE EFFICIENCY METRICS ARE COVERED FROM TWO PERSPECTIVES → COMPUTATION RELATED → MACS

↓  
MEMORY RELATED  
MODEL SIZE ← NUMBER OF PARAMETERS

· THE TOTAL ACTIVATIONS AND THE PEAK ACTIVATIONS

INPUT ACTIVATION SIZE + OUTPUT ACTIVATION SIZE → NN SPECIFICATION  
MEMORY BANDWIDTH OF PROCESSOR → HARDWARE SPECIFICATION

HOW TO CALCULATE THE LATENCY →  $\hat{=}$  MAX ( $T_{\text{COMPUTATION}}$ ,  $T_{\text{MEMORY}}$ ) →  $\hat{=}$   $T_{\text{DATA MOVEMENT OF ACTIVATIONS}} + T_{\text{DATA MOVEMENT OF WEIGHTS}}$  →  $\frac{\text{NUMBER OF OPERATIONS IN NEURAL NETWORKS MODEL}}{\text{NUMBER OF OPERATIONS THAT PROCESSOR CAN PROCESS PER SECOND}}$  →  $\frac{\text{NEURAL NETWORK MODEL SIZE}}{\text{MEMORY BANDWIDTH OF PROCESSOR}}$  → NN SPECIFICATION  
HARDWARE SPECIFICATION → HARDWARE SPECIFICATION

ENERGY CONSUMPTION → MEMORY MOVEMENT IS MUCH MORE EXPENSIVE THAN COMPUTATION! → WE WANT TO REDUCE THE AMOUNT OF DATA MOVEMENT WHEN WE ARE DESIGNING NEURAL NETS.

# OF PARAMETERS → IT'S THE NUMBER OF WEIGHTS FOR A GIVEN NEURAL NETWORK

MODEL SIZE → IT MEASURES THE STORAGE FOR THE WEIGHTS OF A GIVEN NN.

· THE COMMON UNITS FOR THE MODEL SIZE ARE MEGABYTE (MB), KILOBYTE (KB) AND BITS.

· IN GENERAL (IF THE WHOLE NN USES THE SAME DATA TYPE : MODEL SIZE = #PARAMETERS × BIT WIDTH

WE CAN USE 32-BIT FLOATING POINT OR 16 ETC..

# OF ACTIVATIONS → IS THE MEMORY BOTTLENECK IN INFERENCE ON IOT

THE DISTRIBUTION OF THE ACTIVATION IS HIGHLY IMBALANCED → SOME OF THE LAYERS CAN HAVE A DOMINATING ACTIVATION → SO WHAT IS THE LEAST AMOUNT OF MEMORY TO PROCESS THIS NN? THAT IS IMPACTED BY THE PEAK FOR INFERENCE, FOR TRAINING THAT'S IMPACT BY THE SUM.

↑  
FOR TRAINING THE ACTIVATION IS THE BOTTLENECK (NOT # OF PARAMETERS)

MACS → NUMBER OF MULTIPLY-ACCUMULATE OPERATIONS

LAYER

MACS  
(BATCH SIZE = 1)

LINEAR

$C_o \cdot C_i$

CONVENTION

$C_i \cdot K_h \cdot K_w \cdot P_o \cdot W_o \cdot C_o$

GROUPED CONVENTION

$C_i/g \cdot K_h \cdot K_w \cdot P_o \cdot W_o \cdot C_o$

DEPTHWISE CONVENTION

$K_h \cdot K_w \cdot P_o \cdot W_o \cdot C_o$

FLOP → # OF FLOATING POINT OPERATIONS

· A MULTIPLY IS A FLOATING POINT OPERATION AND ADD AN ADD.

· ONE MULTIPLY-ACCUMULATION OPERATION IS TWO FLOATING POINT OPERATIONS.