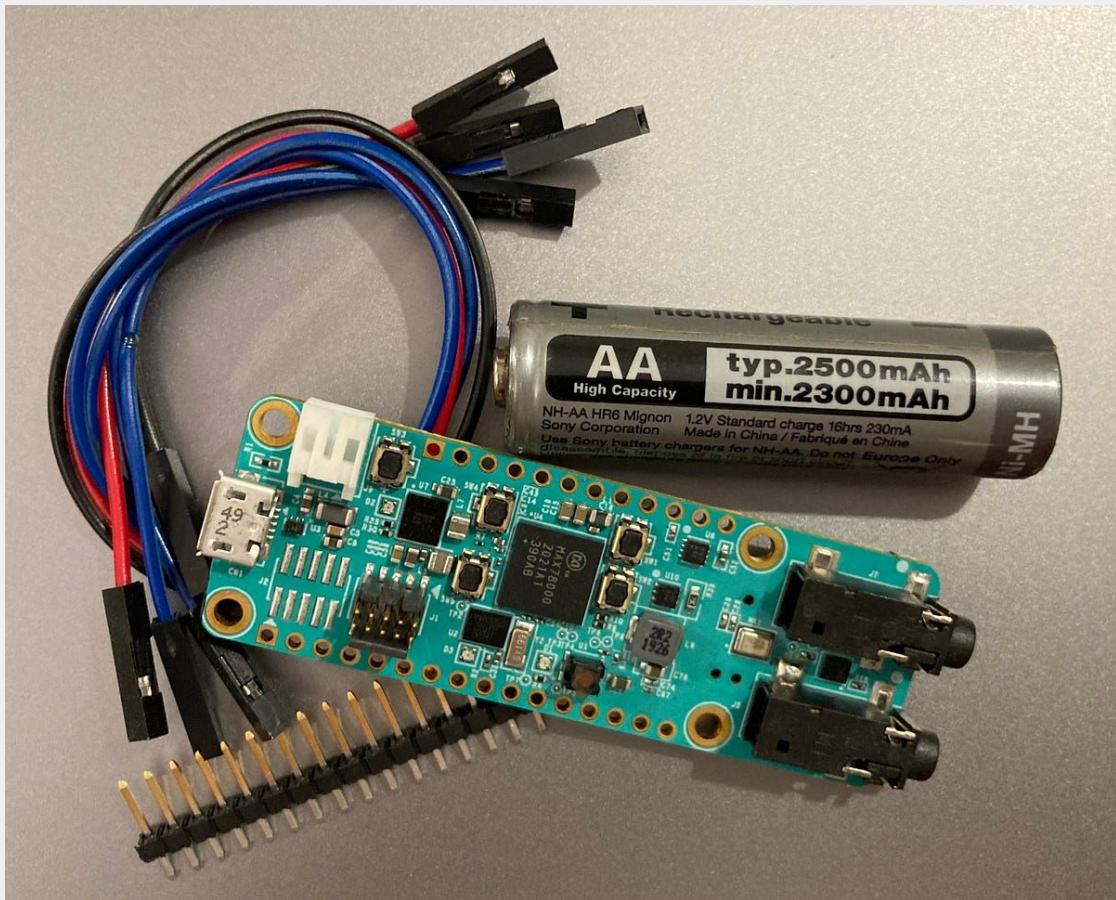


AI on a “Battery”

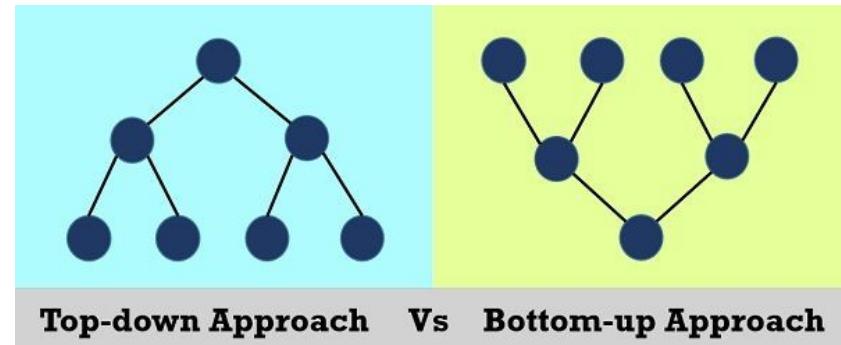
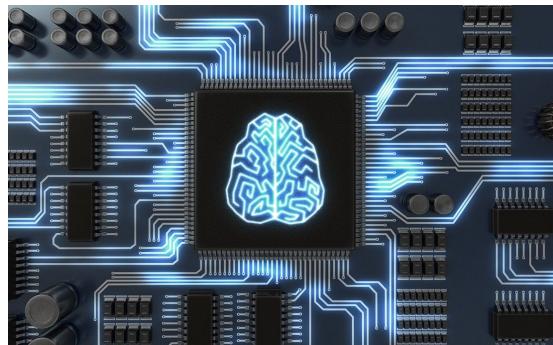


Edge Ai / Ai (on) the Edge

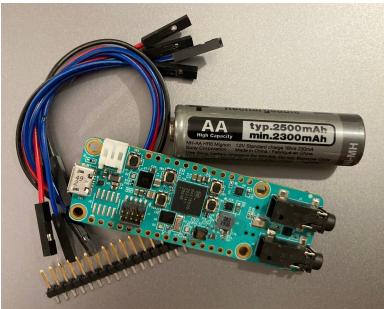
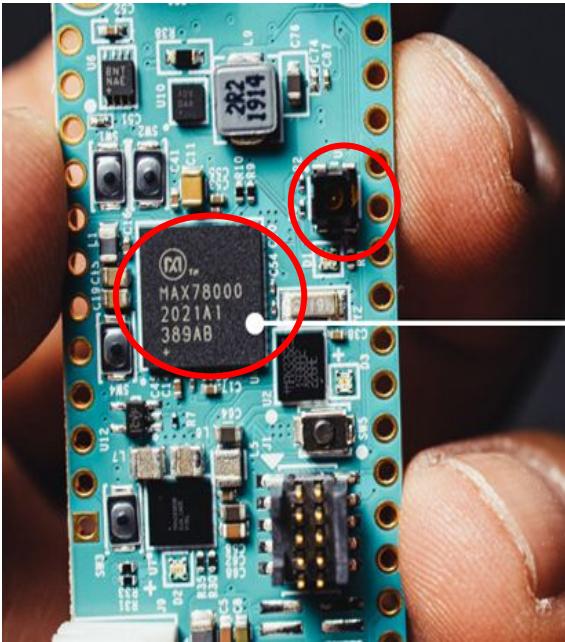
- Traditional CPUs, while versatile, struggled to handle the intensive computations
- specialized hardware accelerators designed to optimize AI processing
 - GPU - ideal for training and running neural networks
 - TPUs (Tensor Processing Units) - custom-designed chips optimized for tensor operations (google)
 - FPGAs - hardware that can be reconfigured to optimize specific tasks
 - NPUs (Neural Processing Units):
coprocessor integrated into various devices designed specifically for AI tasks,
providing high efficiency and performance with minimal power consumption.

ASICs ... dsp ... * NLP

◦



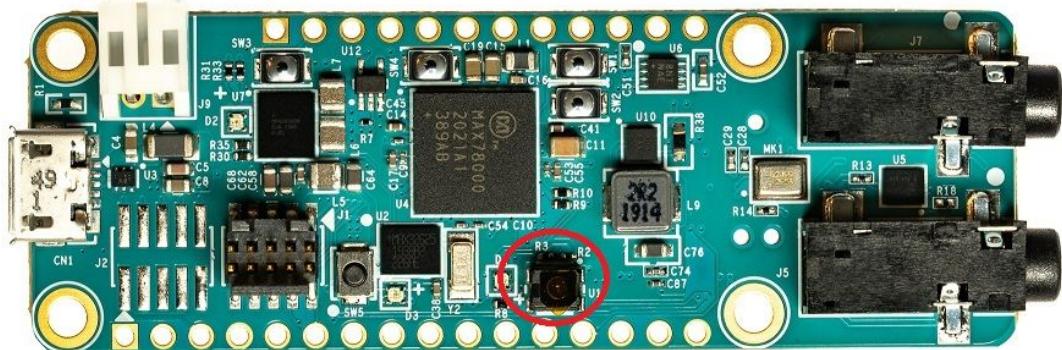
The target



“AI on a Battery”

MAX78000FTHR
Camera module:

Himax HM01B0
320 x 320 pixels

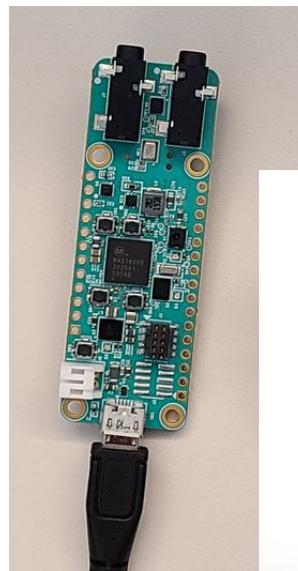


Make sure the protective film on the lens is removed!

- **MCU Core:** ARM Cortex-M4 with FPU.
- **Memory:** **512KB Flash, 128KB SRAM.**
- **Energy Efficiency:** **Ultra-low-power** design, optimized for battery-powered devices.
- **Power Consumption:** Typically around **1 mW** when running inference.
- **AI Performance:** Capable of running AI models like CNNs efficiently on the edge.
 - **TOPS (Tera Operations Per Second):** Approximately 0.1 TOPS (100 GOPS).

blink outdoor 4

(some) application



segmentation

...objective

Python on dev host

```
/data/tinyai/maxCNN/ai8x-train  
schedule-cifar-nas.yaml -model ai85nascifarnet --dataset CIFAR10 --device MAX78000 --batch-size 100 --print-freq 100 --validation-split 0 --use-bias --qat-policy policies/qat_policy_late_cifar.yaml --confusion
```

```
"args": [
```

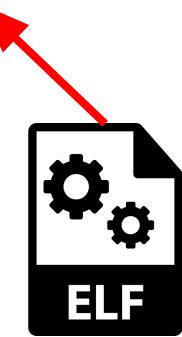
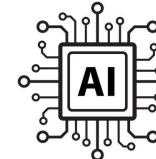
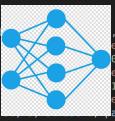
```
--deterministic,  
// "--epochs", "300",  
"--epochs", "1",  
"--optimizer", "Adam",  
"--lr", ".001",  
"--wd", "0",  
"--compress", "policies/schedule-cifar-nas.yaml",  
"--model", "ai85nascifarnet",  
"--dataset", "CIFAR10",  
"--device", "MAX78000",  
"--batch-size", "100",
```



```
--wd 0 --compress policies/s  
--validation-split 0 --use-bias --qat-policy policies/qat_policy_late_cifar.yaml --confusion
```

Python model for ia8x

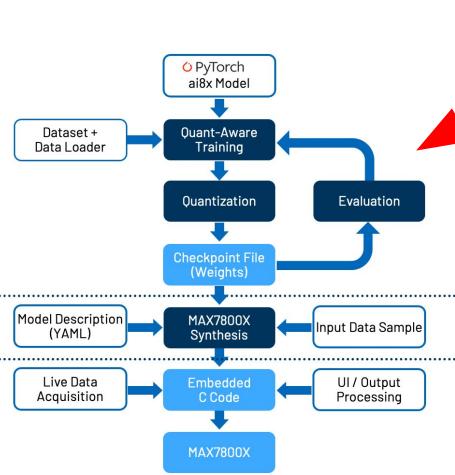
```
class Model(nn.Module):  
    def __init__(self, num_classes=10, num_channels=3, dimensions=(32, 32), **kwargs):  
        super().__init__()  
  
        self.conv1_1 = ai8x.FusedConv2dBNReLU(num_channels, 32, 3, 1, padding=1, bias=bias, **kwargs)  
        self.conv1_2 = ai8x.FusedConv2dBNReLU(32, 32, 1, 1, padding=0, bias=bias, **kwargs)  
        self.conv1_3 = ai8x.FusedConv2dBNReLU(32, 64, 3, 1, 1, 1, padding=1, bias=bias, **kwargs)  
        self.conv2_1 = ai8x.FusedMaxPoolConv2dBNReLU(64, 64, 2, 2, 1, 1, padding=1, bias=bias, batchnorm='NoAffine', **kwargs)  
        self.conv2_2 = ai8x.FusedConv2dBNReLU(32, 64, 1, stride=1, padding=0, bias=bias, batchnorm='NoAffine', **kwargs)  
        self.conv3_1 = ai8x.FusedMaxPoolConv2dBNReLU(64, 128, 3, stride=1, padding=1, bias=bias, batchnorm='NoAffine', **kwargs)  
        self.conv3_2 = ai8x.FusedConv2dBNReLU(128, 128, 1, stride=1, padding=0, bias=bias, batchnorm='NoAffine', **kwargs)  
        self.conv4_1 = ai8x.FusedMaxPoolConv2dBNReLU(128, 64, 3, stride=1, padding=1, bias=bias, batchnorm='NoAffine', **kwargs)
```



C for ARM

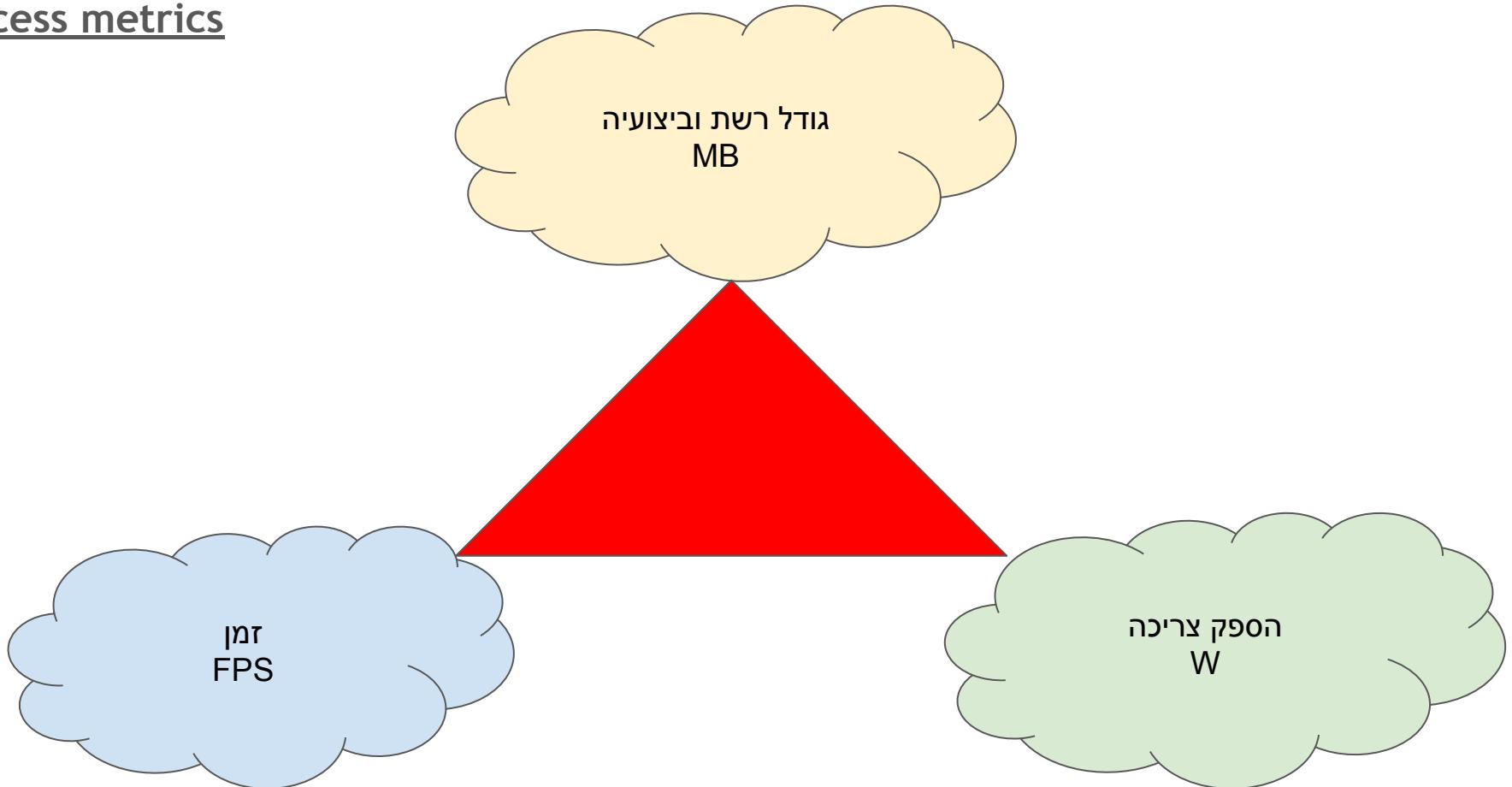
```
#include <ml.h>  
#include <math.h>  
  
float ml_data[CNN_NUM_OUTPUTS];  
float ml_softmax(CNN_NUM_OUTPUTS);  
  
layer(void)  
{  
    if((uint32_t *) ml_data);  
    l7p14_q15((const q31_t *) ml_data, CNN_NUM_OUTPUTS, ml_softmax);  
}
```

```
int main(void)  
{  
    int i;  
    int digs, tens;  
  
    MXC_ICC_Enable(MXC_ICC0); // Enable cache  
  
    // Switch to 100 MHz clock  
    MXC_SYS_Clock_Select(MXC_SYS_CLOCK_IPO);  
    SystemCoreClockUpdate();  
  
    printf("Waiting...\n");  
  
    // DO NOT DELETE THIS LINE:  
    MXC_Delay(SEC(2)); // Let debugger interrupt if needed  
  
    // Enable peripheral, enable CNN interrupt, turn on CNN clock  
    // CNN clock: APB (50 MHz) div 1  
    cnn_enable(MXC_S_GCR_PCLKDIV_CNNCLKSEL_PCLK, MXC_S_GCR_PCLKDIV_CNNCLKDIV_DIV1);
```



```
1_train.md  
2_quant.md  
3_eval.md  
4_random_sample.md  
5_yml.md  
6_loader.md  
7_inference.md  
8_burn.md
```

Success metrics



Starting point:

The screenshot shows a computer screen displaying an arXiv preprint titled "SimpleNet". The page header includes the arXiv logo, the category "cs > Computer Vision and Pattern Recognition", and the submission date "22 Aug 2016 (v1), last revised 27 Apr 2023 (this version, v8)". The main title is "Lets keep it simple, Using simple architectures to outperform deeper and more complex architectures". Below the title, the authors are listed as Seyyed Hossein Hasanpour, Mohammad Rouhani, Mohsen Fayyaz, and Mohammad Sabokrou. The abstract discusses the limitations of deep neural networks like AlexNet, VGGNet, ResNet, and GoogleNet, and introduces SimpleNet, which achieves state-of-the-art results on CIFAR10 while being much simpler and faster. The code section shows the Python implementation of SimpleNet, specifically the AI85SimpleNet class. The code uses the ai8x library for fused convolutional layers. The code is annotated with line numbers from 20 to 42.

```
20 Gorkem Ulkar, 3 years ago | 2
21 class AI85SimpleNet(n
22     """
23     SimpleNet v1 Model
24     """
25     def __init__(
26         self,
27             num_classes=100,
28             num_channels=3,
29             dimensions=(32, 32), # pylint: disable=unused-argument
30             bias=False,
31             **kwargs
32     ):
33         super().__init__()
34
35         self.conv1 = ai8x.FusedConv2dBNReLU(num_channels, 16, 3, stride=1, padding=1, bias=bias,
36                                         **kwargs)
37         self.conv2 = ai8x.FusedConv2dBNReLU(16, 20, 3, stride=1, padding=1, bias=bias, **kwargs)
38         self.conv3 = ai8x.FusedConv2dBNReLU(20, 20, 3, stride=1, padding=1, bias=bias, **kwargs)
39         self.conv4 = ai8x.FusedConv2dBNReLU(20, 20, 3, stride=1, padding=1, bias=bias, **kwargs)
40         self.conv5 = ai8x.FusedMaxPoolConv2dBNReLU(20, 20, 3, pool_size=2, pool_stride=2,
41                                         stride=1, padding=1, bias=bias, **kwargs)
42         self.conv6 = ai8x.FusedConv2dBNReLU(20, 20, 3, stride=1, padding=1, bias=bias, **kwargs)
```

Comments: Added the long-overdue ImageNet results and updated the missed cifar10/100 results from 2018
Subjects: Computer Vision and Pattern Recognition (cs.CV); Neural and Evolutionary Computing (cs.NE)
Cite as: arXiv:1608.06037 [cs.CV]
(or arXiv:1608.06037v8 [cs.CV] for this version)
<https://doi.org/10.48550/arXiv.1608.06037>

Data set : CamVid sXX YY

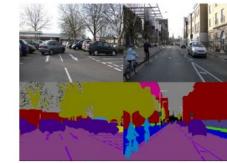
```
train_loader len = 369  
val_loader len = 100  
test_loader len = 232
```

CamVid (Cambridge-driving Labeled Video Database)

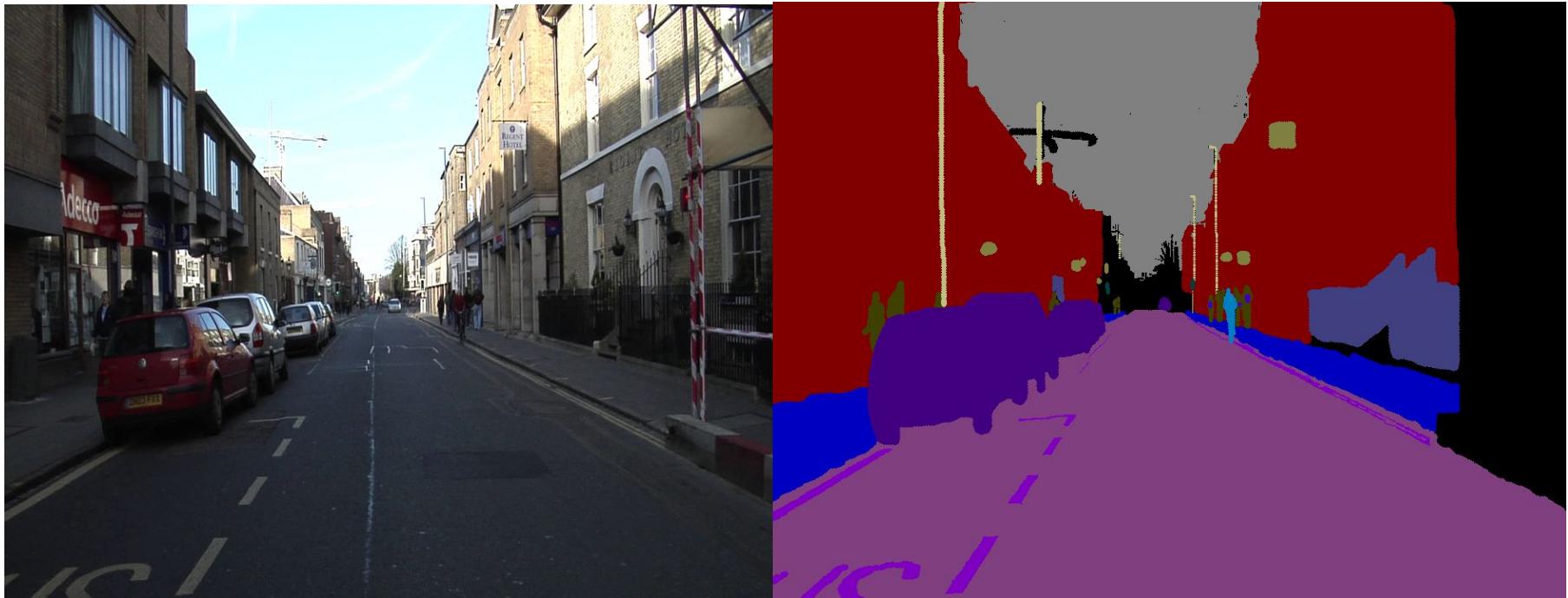
Introduced by Gabriel J. Brostow et al. in [Semantic object classes in video: A high-definition ground truth database](#)

Edit

CamVid (Cambridge-driving Labeled Video Database) is a road/driving scene understanding database which was originally captured as five video sequences with a 960x720 resolution camera mounted on the dashboard of a car. Those sequences were sampled (four of them at 1 fps and one at 15 fps) adding up to 701 frames. Those stills were manually annotated with 32 classes: void, building, wall, tree, vegetation, fence, sidewalk, parking block, column/pole, traffic cone, bridge, sign, miscellaneous text, traffic light, sky, tunnel, archway, road, road shoulder, lane markings (driving), lane markings (non-driving), animal, pedestrian, child, cart luggage, bicyclist, motorcycle, car, SUV/pickup/truck, truck/bus, train, and other moving object



Source: <http://mi.eng.cam.ac.uk/resea...>



Starting point:

A

1. Input Preparation Layers:

- **prep0, prep1, prep2**: These initial layers sequentially transform the input image (3 channels) through a series of convolutional operations:
 - **Conv2d + BatchNorm + ReLU**: Convert the input to higher dimensions (64 -> 64 -> 32 channels).

2. Encoder Path:

- **enc1**: Applies convolution to reduce dimensions (32 to 8 channels).
- **enc2**: MaxPooling followed by convolution (8 to 28 channels).
- **enc3**: Another MaxPooling and convolution step (28 to 56 channels).

3. Bottleneck Layer:

- **bneck**: Central part of the network, combining MaxPooling and convolution (56 to 112 channels).

4. Decoder Path:

- **upconv3**: Upsampling (deconvolution) layer increasing dimensions (112 to 56 channels).
- **dec3**: Combines features from encoder (concatenation) and applies convolution (56 channels).
- **upconv2**: Further upsampling (56 to 28 channels).
- **dec2**: Applies convolution after concatenation (28 channels).
- **upconv1**: Final upsampling layer (28 to 8 channels).
- **dec1**: Final set of convolutional operations after concatenation (16 to 48 channels).
- **dec0**: Final convolution operations (48 to 64 channels).

5. Final Convolution Layers:

- **conv_p1, conv_p2, conv_p3**: Further refine the output with convolutional layers (64 channels each).
- **conv**: Final layer, converting the refined output to the final desired output (64 to 256 channels).

fold

the final number of output channels is calculated as
`num_final_channels = num_classes * fold_ratio * fold_ratio.`

This allows the model to provide a detailed and high-dimensional feature representation, which is beneficial for tasks requiring fine-grained segmentation.

Param size

Pytorch summary

```
)  
    (enc1): Sequential(  
        (0): Conv2d(32, 8, kernel_size=(3, 3), stride=(1, 1), padding=0,  
        (1): BatchNorm2d(8, eps=1e-05, momentum=0.1, affine=True,  
        (2): ReLU(inplace=True)  
)  
    (enc2): Sequential(  
        (0): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1)  
        (1): Conv2d(8, 28, kernel_size=(3, 3), stride=(1, 1), padding=0,  
        (2): BatchNorm2d(28, eps=1e-05, momentum=0.1, affine=True,  
    ...  
Forward/backward pass size (MB): 76.73  
Params size (MB): 1.06  
Estimated Total Size (MB): 77.86
```

Forward/backward pass size (MB):

This is the amount of memory required to store the intermediate activations (outputs of each layer) and gradients during the forward and backward passes of the model. This memory is needed for computing gradients during the backpropagation step.

Params size (MB):

This is the memory required to store the parameters (weights and biases) of the model. This size is computed based on the total number of parameters and their data type (usually float32).

Estimated Total Size:

total memory needed to train the model

```
# Calculate number of parameters  
total_params, trainable_params = count_parameters(model)  
print(f"Trainable parameters: {trainable_params}")  
  
def dtype_to_bytes(dtype):  
    dtype_bytes = {  
        torch.float32: 4,  
        torch.float64: 8,  
        torch.float16: 2,  
        torch.int32: 4,  
        torch.int64: 8  
    }  
    # Add other data types if needed  
}
```

```
check_unique_param_sizes(model)
```

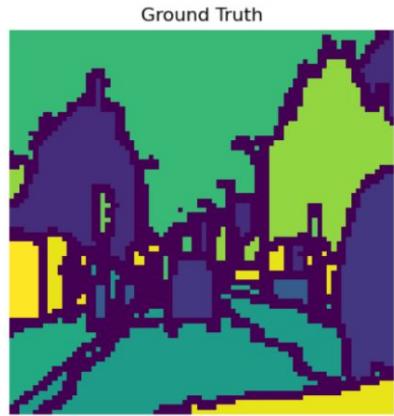
Size: 4 bytes, Count: 66

```
size_of_each_param_in_bytes=4  
total_size_in_bytes=total_params * size_of_each_param_in_bytes  
total_size_in_megabytes = total_size_in_bytes / (1024*1024)  
  
print(f"total_size_in_megabytes = {total_size_in_megabytes}")  
✓ 0.0s  
total_size_in_megabytes = 1.0598602294921875
```

1try:

c33

```
{  
    'name': 'CamVid_s352_c33',  
    'input': (48, 88, 88),  
    'output': ('None', 'Animal', 'Archway', 'Bicyclist', 'Bridge', 'Building', 'Car',  
              'CartLuggagePram', 'Child', 'Column_Pole', 'Fence', 'LaneMkgsDriv',  
              'LaneMkgsNonDriv', 'Misc_Text', 'MotorcycleScooter', 'OtherMoving',  
              'ParkingBlock', 'Pedestrian', 'Road', 'RoadShoulder', 'Sidewalk', 'SignSymbol',  
              'Sky', 'SUVPickupTruck', 'TrafficCone', 'TrafficLight', 'Train', 'Tree',  
              'Truck_Bus', 'Tunnel', 'VegetationMisc', 'Void', 'Wall'),  
    'loader': camvid_get_datasets_s352_c33,  
    'fold_ratio': 4,  
},
```



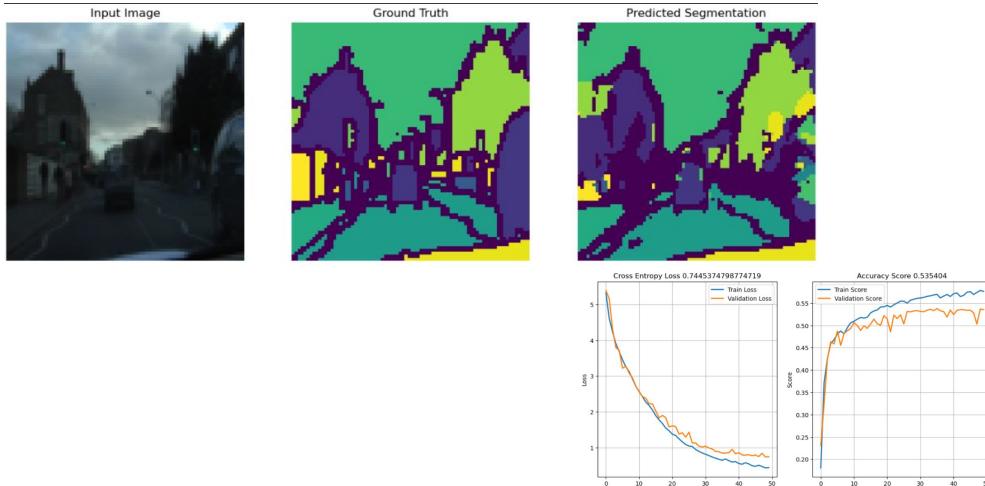
c3

```
{  
    'name': 'CamVid_s80_c3', # 3 classes  
    'input': (3, 80, 80), # 23  
    'output': (0, 1, 2), # 24  
    'weight': (1, 1, 1), # 25  
    'loader': camvid_get_datasets_s80_c3, # 26  
    'filter_label': [5, 21, 17], # car sky road  
},  
    # 27  
    # 28  
    # 29  
    # 30  
    # 31  
    # 32  
    # 33  
    # 34  
    # 35  
    # 36  
    # 37  
    # 38  
    # 39  
    # 40  
    # 41  
    # 42  
    # 43  
    # 44  
    # 45  
    # 46  
    # 47  
    # 48  
    # 49  
    # 50  
    # 51  
    # 52  
    # 53  
    # 54  
    # 55  
    # 56  
    # 57  
    # 58  
    # 59  
    # 60  
    # 61  
    # 62  
    # 63  
    # 64  
    # 65  
    # 66  
    # 67  
    # 68  
    # 69  
    # 70  
    # 71  
    # 72  
    # 73  
    # 74  
    # 75  
    # 76  
    # 77  
    # 78  
    # 79  
    # 80  
    # 81  
    # 82  
    # 83  
    # 84  
    # 85  
    # 86  
    # 87  
    # 88  
    # 89  
    # 90  
    # 91  
    # 92  
    # 93  
    # 94  
    # 95  
    # 96  
    # 97  
    # 98  
    # 99  
    # 100  
    # 101  
    # 102  
    # 103  
    # 104  
    # 105  
    # 106  
    # 107  
    # 108  
    # 109  
    # 110  
    # 111  
    # 112  
    # 113  
    # 114  
    # 115  
    # 116  
    # 117  
    # 118  
    # 119  
    # 120  
    # 121  
    # 122  
    # 123  
    # 124  
    # 125  
    # 126  
    # 127  
    # 128  
    # 129  
    # 130  
    # 131  
    # 132  
    # 133  
    # 134  
    # 135  
    # 136  
    # 137  
    # 138  
    # 139  
    # 140  
    # 141  
    # 142  
    # 143  
    # 144  
    # 145  
    # 146  
    # 147  
    # 148  
    # 149  
    # 150  
    # 151  
    # 152  
    # 153  
    # 154  
    # 155  
    # 156  
    # 157  
    # 158  
    # 159  
    # 160  
    # 161  
    # 162  
    # 163  
    # 164  
    # 165  
    # 166  
    # 167  
    # 168  
    # 169  
    # 170  
    # 171  
    # 172  
    # 173  
    # 174  
    # 175  
    # 176  
    # 177  
    # 178  
    # 179  
    # 180  
    # 181  
    # 182  
    # 183  
    # 184  
    # 185  
    # 186  
    # 187  
    # 188  
    # 189  
    # 190  
    # 191  
    # 192  
    # 193  
    # 194  
    # 195  
    # 196  
    # 197  
    # 198  
    # 199  
    # 200  
    # 201  
    # 202  
    # 203  
    # 204  
    # 205  
    # 206  
    # 207  
    # 208  
    # 209  
    # 210  
    # 211  
    # 212  
    # 213  
    # 214  
    # 215  
    # 216  
    # 217  
    # 218  
    # 219  
    # 220  
    # 221  
    # 222  
    # 223  
    # 224  
    # 225  
    # 226  
    # 227  
    # 228  
    # 229  
    # 230  
    # 231  
    # 232  
    # 233  
    # 234  
    # 235  
    # 236  
    # 237  
    # 238  
    # 239  
    # 240  
    # 241  
    # 242  
    # 243  
    # 244  
    # 245  
    # 246  
    # 247  
    # 248  
    # 249  
    # 250  
    # 251  
    # 252  
    # 253  
    # 254  
    # 255  
    # 256  
    # 257  
    # 258  
    # 259  
    # 260  
    # 261  
    # 262  
    # 263  
    # 264  
    # 265  
    # 266  
    # 267  
    # 268  
    # 269  
    # 270  
    # 271  
    # 272  
    # 273  
    # 274  
    # 275  
    # 276  
    # 277  
    # 278  
    # 279  
    # 280  
    # 281  
    # 282  
    # 283  
    # 284  
    # 285  
    # 286  
    # 287  
    # 288  
    # 289  
    # 290  
    # 291  
    # 292  
    # 293  
    # 294  
    # 295  
    # 296  
    # 297  
    # 298  
    # 299  
    # 300  
    # 301  
    # 302  
    # 303  
    # 304  
    # 305  
    # 306  
    # 307  
    # 308  
    # 309  
    # 310  
    # 311  
    # 312  
    # 313  
    # 314  
    # 315  
    # 316  
    # 317  
    # 318  
    # 319  
    # 320  
    # 321  
    # 322  
    # 323  
    # 324  
    # 325  
    # 326  
    # 327  
    # 328  
    # 329  
    # 330  
    # 331  
    # 332  
    # 333  
    # 334  
    # 335  
    # 336  
    # 337  
    # 338  
    # 339  
    # 340  
    # 341  
    # 342  
    # 343  
    # 344  
    # 345  
    # 346  
    # 347  
    # 348  
    # 349  
    # 350  
    # 351  
    # 352  
    # 353  
    # 354  
    # 355  
    # 356  
    # 357  
    # 358  
    # 359  
    # 360  
    # 361  
    # 362  
    # 363  
    # 364  
    # 365  
    # 366  
    # 367  
    # 368  
    # 369  
    # 370  
    # 371  
    # 372  
    # 373  
    # 374  
    # 375  
    # 376  
    # 377  
    # 378  
    # 379  
    # 380  
    # 381  
    # 382  
    # 383  
    # 384  
    # 385  
    # 386  
    # 387  
    # 388  
    # 389  
    # 390  
    # 391  
    # 392  
    # 393  
    # 394  
    # 395  
    # 396  
    # 397  
    # 398  
    # 399  
    # 400  
    # 401  
    # 402  
    # 403  
    # 404  
    # 405  
    # 406  
    # 407  
    # 408  
    # 409  
    # 410  
    # 411  
    # 412  
    # 413  
    # 414  
    # 415  
    # 416  
    # 417  
    # 418  
    # 419  
    # 420  
    # 421  
    # 422  
    # 423  
    # 424  
    # 425  
    # 426  
    # 427  
    # 428  
    # 429  
    # 430  
    # 431  
    # 432  
    # 433  
    # 434  
    # 435  
    # 436  
    # 437  
    # 438  
    # 439  
    # 440  
    # 441  
    # 442  
    # 443  
    # 444  
    # 445  
    # 446  
    # 447  
    # 448  
    # 449  
    # 450  
    # 451  
    # 452  
    # 453  
    # 454  
    # 455  
    # 456  
    # 457  
    # 458  
    # 459  
    # 460  
    # 461  
    # 462  
    # 463  
    # 464  
    # 465  
    # 466  
    # 467  
    # 468  
    # 469  
    # 470  
    # 471  
    # 472  
    # 473  
    # 474  
    # 475  
    # 476  
    # 477  
    # 478  
    # 479  
    # 480  
    # 481  
    # 482  
    # 483  
    # 484  
    # 485  
    # 486  
    # 487  
    # 488  
    # 489  
    # 490  
    # 491  
    # 492  
    # 493  
    # 494  
    # 495  
    # 496  
    # 497  
    # 498  
    # 499  
    # 500  
    # 501  
    # 502  
    # 503  
    # 504  
    # 505  
    # 506  
    # 507  
    # 508  
    # 509  
    # 510  
    # 511  
    # 512  
    # 513  
    # 514  
    # 515  
    # 516  
    # 517  
    # 518  
    # 519  
    # 520  
    # 521  
    # 522  
    # 523  
    # 524  
    # 525  
    # 526  
    # 527  
    # 528  
    # 529  
    # 530  
    # 531  
    # 532  
    # 533  
    # 534  
    # 535  
    # 536  
    # 537  
    # 538  
    # 539  
    # 540  
    # 541  
    # 542  
    # 543  
    # 544  
    # 545  
    # 546  
    # 547  
    # 548  
    # 549  
    # 550  
    # 551  
    # 552  
    # 553  
    # 554  
    # 555  
    # 556  
    # 557  
    # 558  
    # 559  
    # 560  
    # 561  
    # 562  
    # 563  
    # 564  
    # 565  
    # 566  
    # 567  
    # 568  
    # 569  
    # 570  
    # 571  
    # 572  
    # 573  
    # 574  
    # 575  
    # 576  
    # 577  
    # 578  
    # 579  
    # 580  
    # 581  
    # 582  
    # 583  
    # 584  
    # 585  
    # 586  
    # 587  
    # 588  
    # 589  
    # 590  
    # 591  
    # 592  
    # 593  
    # 594  
    # 595  
    # 596  
    # 597  
    # 598  
    # 599  
    # 600  
    # 601  
    # 602  
    # 603  
    # 604  
    # 605  
    # 606  
    # 607  
    # 608  
    # 609  
    # 610  
    # 611  
    # 612  
    # 613  
    # 614  
    # 615  
    # 616  
    # 617  
    # 618  
    # 619  
    # 620  
    # 621  
    # 622  
    # 623  
    # 624  
    # 625  
    # 626  
    # 627  
    # 628  
    # 629  
    # 630  
    # 631  
    # 632  
    # 633  
    # 634  
    # 635  
    # 636  
    # 637  
    # 638  
    # 639  
    # 640  
    # 641  
    # 642  
    # 643  
    # 644  
    # 645  
    # 646  
    # 647  
    # 648  
    # 649  
    # 650  
    # 651  
    # 652  
    # 653  
    # 654  
    # 655  
    # 656  
    # 657  
    # 658  
    # 659  
    # 660  
    # 661  
    # 662  
    # 663  
    # 664  
    # 665  
    # 666  
    # 667  
    # 668  
    # 669  
    # 670  
    # 671  
    # 672  
    # 673  
    # 674  
    # 675  
    # 676  
    # 677  
    # 678  
    # 679  
    # 680  
    # 681  
    # 682  
    # 683  
    # 684  
    # 685  
    # 686  
    # 687  
    # 688  
    # 689  
    # 690  
    # 691  
    # 692  
    # 693  
    # 694  
    # 695  
    # 696  
    # 697  
    # 698  
    # 699  
    # 700  
    # 701  
    # 702  
    # 703  
    # 704  
    # 705  
    # 706  
    # 707  
    # 708  
    # 709  
    # 710  
    # 711  
    # 712  
    # 713  
    # 714  
    # 715  
    # 716  
    # 717  
    # 718  
    # 719  
    # 720  
    # 721  
    # 722  
    # 723  
    # 724  
    # 725  
    # 726  
    # 727  
    # 728  
    # 729  
    # 730  
    # 731  
    # 732  
    # 733  
    # 734  
    # 735  
    # 736  
    # 737  
    # 738  
    # 739  
    # 740  
    # 741  
    # 742  
    # 743  
    # 744  
    # 745  
    # 746  
    # 747  
    # 748  
    # 749  
    # 750  
    # 751  
    # 752  
    # 753  
    # 754  
    # 755  
    # 756  
    # 757  
    # 758  
    # 759  
    # 760  
    # 761  
    # 762  
    # 763  
    # 764  
    # 765  
    # 766  
    # 767  
    # 768  
    # 769  
    # 770  
    # 771  
    # 772  
    # 773  
    # 774  
    # 775  
    # 776  
    # 777  
    # 778  
    # 779  
    # 780  
    # 781  
    # 782  
    # 783  
    # 784  
    # 785  
    # 786  
    # 787  
    # 788  
    # 789  
    # 790  
    # 791  
    # 792  
    # 793  
    # 794  
    # 795  
    # 796  
    # 797  
    # 798  
    # 799  
    # 800  
    # 801  
    # 802  
    # 803  
    # 804  
    # 805  
    # 806  
    # 807  
    # 808  
    # 809  
    # 810  
    # 811  
    # 812  
    # 813  
    # 814  
    # 815  
    # 816  
    # 817  
    # 818  
    # 819  
    # 820  
    # 821  
    # 822  
    # 823  
    # 824  
    # 825  
    # 826  
    # 827  
    # 828  
    # 829  
    # 830  
    # 831  
    # 832  
    # 833  
    # 834  
    # 835  
    # 836  
    # 837  
    # 838  
    # 839  
    # 840  
    # 841  
    # 842  
    # 843  
    # 844  
    # 845  
    # 846  
    # 847  
    # 848  
    # 849  
    # 850  
    # 851  
    # 852  
    # 853  
    # 854  
    # 855  
    # 856  
    # 857  
    # 858  
    # 859  
    # 860  
    # 861  
    # 862  
    # 863  
    # 864  
    # 865  
    # 866  
    # 867  
    # 868  
    # 869  
    # 870  
    # 871  
    # 872  
    # 873  
    # 874  
    # 875  
    # 876  
    # 877  
    # 878  
    # 879  
    # 880  
    # 881  
    # 882  
    # 883  
    # 884  
    # 885  
    # 886  
    # 887  
    # 888  
    # 889  
    # 890  
    # 891  
    # 892  
    # 893  
    # 894  
    # 895  
    # 896  
    # 897  
    # 898  
    # 899  
    # 900  
    # 901  
    # 902  
    # 903  
    # 904  
    # 905  
    # 906  
    # 907  
    # 908  
    # 909  
    # 910  
    # 911  
    # 912  
    # 913  
    # 914  
    # 915  
    # 916  
    # 917  
    # 918  
    # 919  
    # 920  
    # 921  
    # 922  
    # 923  
    # 924  
    # 925  
    # 926  
    # 927  
    # 928  
    # 929  
    # 930  
    # 931  
    # 932  
    # 933  
    # 934  
    # 935  
    # 936  
    # 937  
    # 938  
    # 939  
    # 940  
    # 941  
    # 942  
    # 943  
    # 944  
    # 945  
    # 946  
    # 947  
    # 948  
    # 949  
    # 950  
    # 951  
    # 952  
    # 953  
    # 954  
    # 955  
    # 956  
    # 957  
    # 958  
    # 959  
    # 960  
    # 961  
    # 962  
    # 963  
    # 964  
    # 965  
    # 966  
    # 967  
    # 968  
    # 969  
    # 970  
    # 971  
    # 972  
    # 973  
    # 974  
    # 975  
    # 976  
    # 977  
    # 978  
    # 979  
    # 980  
    # 981  
    # 982  
    # 983  
    # 984  
    # 985  
    # 986  
    # 987  
    # 988  
    # 989  
    # 990  
    # 991  
    # 992  
    # 993  
    # 994  
    # 995  
    # 996  
    # 997  
    # 998  
    # 999  
    # 1000  
    # 1001  
    # 1002  
    # 1003  
    # 1004  
    # 1005  
    # 1006  
    # 1007  
    # 1008  
    # 1009  
    # 1010  
    # 1011  
    # 1012  
    # 1013  
    # 1014  
    # 1015  
    # 1016  
    # 1017  
    # 1018  
    # 1019  
    # 1020  
    # 1021  
    # 1022  
    # 1023  
    # 1024  
    # 1025  
    # 1026  
    # 1027  
    # 1028  
    # 1029  
    # 1030  
    # 1031  
    # 1032  
    # 1033  
    # 1034  
    # 1035  
    # 1036  
    # 1037  
    # 1038  
    # 1039  
    # 1040  
    # 1041  
    # 1042  
    # 1043  
    # 1044  
    # 1045  
    # 1046  
    # 1047  
    # 1048  
    # 1049  
    # 1050  
    # 1051  
    # 1052  
    # 1053  
    # 1054  
    # 1055  
    # 1056  
    # 1057  
    # 1058  
    # 1059  
    # 1060  
    # 1061  
    # 1062  
    # 1063  
    # 1064  
    # 1065  
    # 1066  
    # 1067  
    # 1068  
    # 1069  
    # 1070  
    # 1071  
    # 1072  
    # 1073  
    # 1074  
    # 1075  
    # 1076  
    # 1077  
    # 1078  
    # 1079  
    # 1080  
    # 1081  
    # 1082  
    # 1083  
    # 1084  
    # 1085  
    # 1086  
    # 1087  
    # 1088  
    # 1089  
    # 1090  
    # 1091  
    # 1092  
    # 1093  
    # 1094  
    # 1095  
    # 1096  
    # 1097  
    # 1098  
    # 1099  
    # 1100  
    # 1101  
    # 1102  
    # 1103  
    # 1104  
    # 1105  
    # 1106  
    # 1107  
    # 1108  
    # 1109  
    # 1110  
    # 1111  
    # 1112  
    # 1113  
    # 1114  
    # 1115  
    # 1116  
    # 1117  
    # 1118  
    # 1119  
    # 1120  
    # 1121  
    # 1122  
    # 1123  
    # 1124  
    # 1125  
    # 1126  
    # 1127  
    # 1128  
    # 1129  
    # 1130  
    # 1131  
    # 1132  
    # 1133  
    # 1134  
    # 1135  
    # 1136  
    # 1137  
    # 1138  
    # 1139  
    # 1140  
    # 1141  
    # 1142  
    # 1143  
    # 1144  
    # 1145  
    # 1146  
    # 1147  
    # 1148  
    # 1149  
    # 1150  
    # 1151  
    # 1152  
    # 1153  
    # 1154  
    # 1155  
    # 1156  
    # 1157  
    # 1158  
    # 1159  
    # 1160  
    # 1161  
    # 1162  
    # 1163  
    # 1164  
    # 1165  
    # 1166  
    # 1167  
    # 1168  
    # 1169  
    # 1170  
    # 1171  
    # 1172  
    # 1173  
    # 1174  
    # 1175  
    # 1176  
    # 1177  
    # 1178  
    # 1179  
    # 1180  
    # 1181  
    # 1182  
    # 1183  
    # 1184  
    # 1185  
    # 1186  
    # 1187  
    # 1188  
    # 1189  
    # 1190  
    # 1191  
    # 1192  
    # 1193  
    # 1194  
    # 1195  
    # 1196  
    # 1197  
    # 1198  
    # 1199  
    # 1200  
    # 1201  
    # 1202  
    # 1203  
    # 1204  
    # 1205  
    # 1206  
    # 1207  
    # 1208  
    # 1209  
    # 1210  
    # 1211  
    # 1212  
    # 1213  
    # 1214  
    # 1215  
    # 1216  
    # 1217  
    # 1218  
    # 1219  
    # 1220  
    # 1221  
    # 1222  
    # 1223  
    # 1224  
    # 1225  
    # 1226  
    # 1227  
    # 1228  
    # 1229  
    # 1230  
    # 1231  
    # 1232  
    # 1233  
    # 1234  
    # 1235  
    # 1236  
    # 1237  
    # 1238  
    # 1239  
    # 1240  
    # 1241  
    # 1242  
    # 1243  
    # 1244  
    # 1245  
    # 1246  
    # 1247  
    # 1248  
    # 1249  
    # 1250  
    # 1251  
    # 1252  
    # 1253  
    # 1254  
    # 1255  
    # 1256  
    # 1257  
    # 1258  
    # 1259  
    # 1260  
    # 1261  
    # 1262  
    # 1263  
    # 1264  
    # 1265  
    # 1266  
    # 1267  
    # 1268  
    # 1269  
    # 1270  
    # 1271  
    # 1272  
    # 1273  
    # 1274  
    # 1275  
    # 1276  
    # 1277  
    # 1278  
    # 1279  
    # 1280  
    # 1281  
    # 1282  
    # 1283  
    # 1284  
    # 1285  
    # 1286  
    # 1287  
    # 1288  
    # 1289  
    # 1290  
    # 1291  
    # 1292  
    # 1293  
    # 1294  
    # 1295  
    # 1296  
    # 1297  
    # 1298  
    # 1299  
    # 1300  
    # 1301  
    # 1302  
    # 1303  
    # 1304  
    # 1305  
    # 1306  
    # 1307  
    # 1308  
    # 1309  
    # 1310  
    # 1311  
    # 1312  
    # 1313  
    # 1314  
    # 1315  
    # 1316  
    # 1317  
    # 1318  
    # 1319  
    # 1320  
    # 1321  
    # 1322  
    # 1323  
    # 1324  
    # 1325  
    # 1326  
    # 1327  
    # 1328  
    # 1329  
    # 1330  
    # 1331  
    # 1332  
    # 1333  
    # 1334  
    # 1335  
    # 1336  
    # 1337  
    # 1338  
    # 1339  
    # 1340  
    # 1341  
    # 1342  
    # 1343  
    # 1344  
    # 1345  
    # 1346  
    # 1347  
    # 1348  
    # 1349  
    # 1350  
    # 1351  
    # 1352  
    # 1353  
    # 1354  
    # 1355  
    # 1356  
    # 1357  
    # 1358  
    # 1359  
    # 1360  
    # 1361  
    # 1362  
    # 1363  
    # 1364  
    # 1365  
    # 1366  
    # 1367  
    # 1368  
    # 1369  
    # 1370  
    # 1371  
    # 1372  
    # 1373  
    # 1374  
    # 1375  
    # 1376  
    # 1377  
    # 1378  
    # 1379  
    # 1380  
    # 1381  
    # 1382  
    # 1383  
    # 1384  
    # 1385  
    # 1386  
    # 1387  
    # 1388  
    # 1389  
    # 1390  
    # 1391  
    # 1392  
    # 1393  
    # 1394  
    # 1395  
    # 1396  
    # 1397  
    # 1398  
    # 1399  
    # 1400  
    # 1401  
    # 1402  
    # 1403  
    # 1404  
    # 1405  
    # 1406  
    # 1407  
    # 1408  
    # 1409  
    # 1410  
    # 1411  
    # 1412  
    # 1413  
    # 1414  
    # 1415  
    # 1416  
    # 1417  
    # 1418  
    # 1419  
    # 1420  
    # 1421  
    # 1422  
    # 1423  
    # 1424  
    # 1425  
    # 1426  
    # 1427  
    # 1428  
    # 1429  
    # 1430  
    # 1431  
    # 1432  
    # 1433  
    # 1434  
    # 1435  
    # 1436  
    # 1437  
    # 1438  
    # 1439  
    # 1440  
    # 1441  
    # 1442  
    # 144
```

1try:

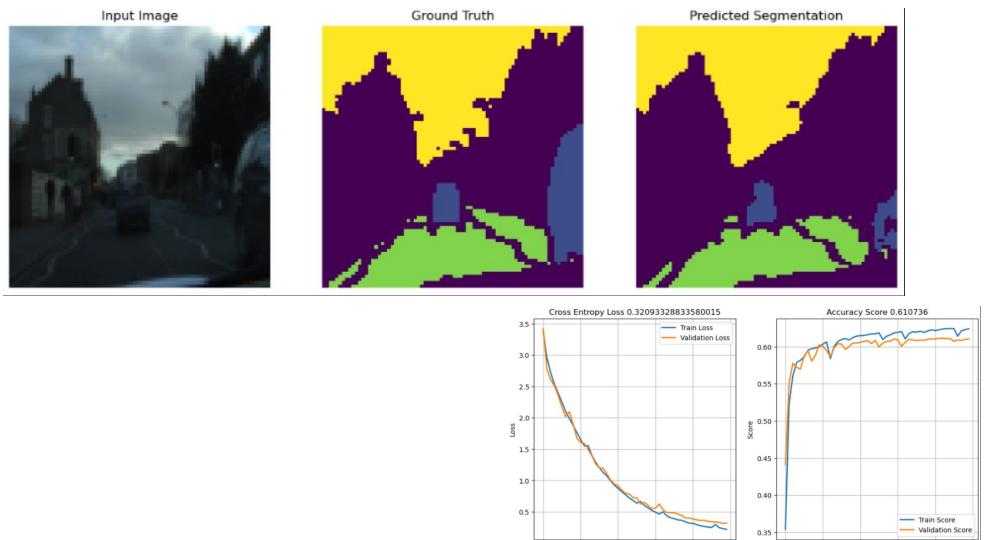
UNetLarge CamVid_80_c33

```
trainning   : EPOCHS_TO_RUN=50 ; batch=8
model       : file=./unet_model_large_x80_c33_e50_A.pth
input data  : train_len = 369 ; test_len=100 ; val_len=232
img data    : IMAGE_SIZE=(80, 80)
model       : NUM_CLASS=33 ; NUM_CHANN=3 ; FOLD_RATIO=4
score        : val_score=0.535404 ; loss_score 0.7445374798774719
Trainable parameters: 308924
parameters_size_in_megabytes : 1.1784515380859375
```



UNetLarge CamVid_80_c33

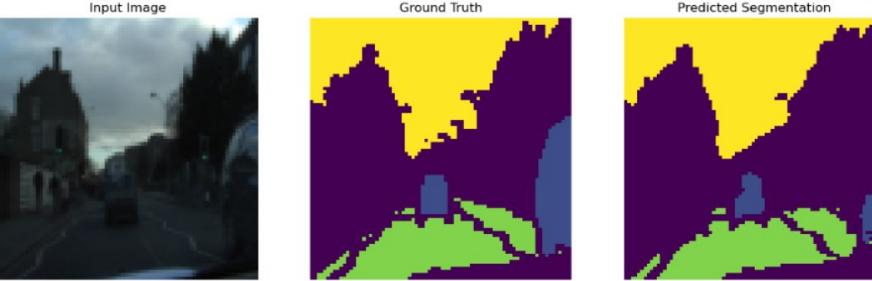
```
trainning   : EPOCHS_TO_RUN=50 ; batch=8
model       : file=./unet_model_large_x80_c3_e50_B.pth
input data  : train_len = 369 ; test_len=100 ; val_len=232
img data    : IMAGE_SIZE=(80, 80)
model       : NUM_CLASS=4 ; NUM_CHANN=3 ; FOLD_RATIO=4
score        : val_score=0.610736 ; loss_score 0.32093328833580015
Trainable parameters: 277836
parameters_size_in_megabytes : 1.0598602294921875
```



2try:

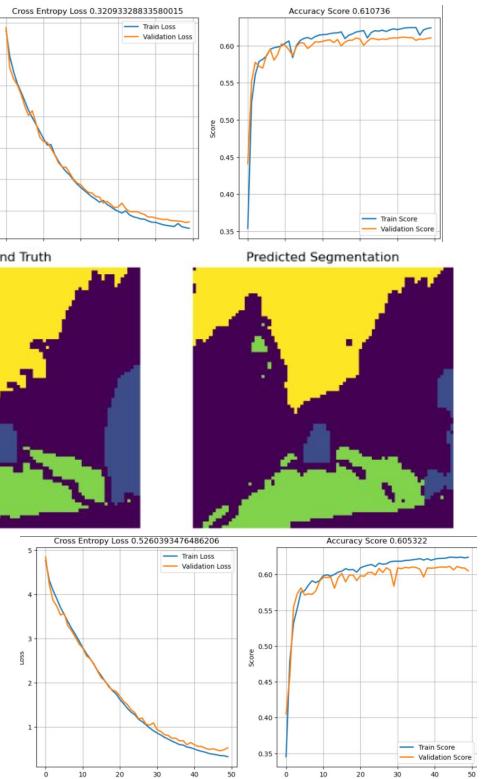
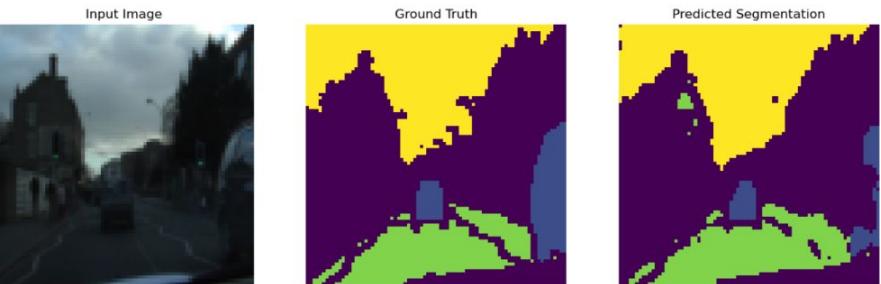
UNetLarge CamVid_80_c33

```
trainning : EPOCHS_TO_RUN=50 ; batch=8
model    : file=./unet_model_large_x80_c3_e50_B.pth
input data : train_len = 369 ; test_len=100 ; val_len=232
img data  : IMAGE_SIZE=(80, 80)
model    : NUM_CLASS=4 ; NUM_CHANN=3 ; FOLD_RATIO=4
score    : val_score=0.610736 ; loss_score 0.32093328833580015
Trainable parameters: 277836
parameters_size_in_megabytes : 1.0598602294921875
```



UNetLarge CamVid_80_c3

```
trainning : EPOCHS_TO_RUN=50 ; batch=8
model    : file=./unet_model_large_x80_c3_e50_C.pth
input data : train_len = 369 ; test_len=100 ; val_len=232
img data  : IMAGE_SIZE=(80, 80)
model    : NUM_CLASS=4 ; NUM_CHANN=3 ; FOLD_RATIO=8
score    : val_score=0.605322 ; loss_score 0.5260393476486206
Trainable parameters: 290700
parameters_size_in_megabytes : 1.1089324951171875
```



1. Initial Layers (prep0, prep1, prep2):

- These layers prepare the input image by applying a series of depthwise separable convolutions, batch normalization, and ReLU activations. They transform the input channels (3 for RGB) to intermediate feature representations (3 -> 32 -> 32 -> 16 channels).

2. Encoder Path (enc1, enc2, enc3):

- The encoder progressively reduces the spatial dimensions while increasing the feature depth:
 - enc1:** Converts 16 channels to 8 through depthwise separable convolutions.
 - enc2:** Downsamples (via max pooling) and then increases features from 8 to 16 channels.
 - enc3:** Further downsampling and feature increase from 16 to 32 channels.

3. Bottleneck Layer (bneck):

- Acts as the central layer, reducing spatial dimensions further while capturing complex features (32 to 64 channels).

4. Decoder Path (upconv3, dec3, upconv2, dec2, upconv1, dec1, dec0):

- The decoder upsamples the feature maps and combines them with corresponding encoder features through concatenation:
 - upconv3 & dec3:** Upsamples from 64 to 32 channels and combines with encoder features.
 - upconv2 & dec2:** Further upsamples and processes features from 32 to 16 channels.
 - upconv1 & dec1:** Final upsampling steps, combining and processing features down to 8 channels and then up to 16 channels.
 - dec0:** Final set of convolutional operations, increasing features from 16 to 32 channels.

5. Final Convolution Layers (conv_p1, conv_p2, conv_p3, conv):

- These layers refine the output further and prepare the final feature map:
 - The final convolution layer outputs 64 channels, converting the processed features to the final desired output.

The number of output channels in the final layer is fixed at 64. provide a detailed feature space, allowing the network to capture intricate details necessary for accurate image segmentation.

diffs

```
[20]    print_model_summary(model, input_size=(3, IMAGE_SIZE[0], IMAGE_SIZE[1]))  
    ✓ 0.0s  
...  
UNetLarge(  
    (prep0): Sequential(  
        (0): Conv2d(3, 64, kernel_size=(1, 1), stride=(1, 1))  
        (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (2): ReLU(inplace=True)  
    )  
    (prep1): Sequential(  
        (0): Conv2d(64, 64, kernel_size=(1, 1), stride=(1, 1))  
        (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (2): ReLU(inplace=True)  
    )  
    (prep2): Sequential(  
        (0): Conv2d(64, 32, kernel_size=(1, 1), stride=(1, 1))  
        (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (2): ReLU(inplace=True)  
    )  
    (enc1): Sequential(  
        (0): Conv2d(32, 8, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (1): BatchNorm2d(8, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (2): ReLU(inplace=True)  
    )  
    (enc2): Sequential(  
        (0): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
        (1): Conv2d(8, 28, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
        (2): BatchNorm2d(28, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (3): ReLU(inplace=True)  
    )  
Python
```

```
[21]    summary(model, input_size=input_size)  
    ✓ 0.0s  
...  
UNetSmall(  
    (prep0): Sequential(  
        (0): Conv2d(3, 3, kernel_size=(1, 1), stride=(1, 1), groups=3)  
        (1): BatchNorm2d(3, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (2): ReLU(inplace=True)  
        (3): Conv2d(3, 32, kernel_size=(1, 1), stride=(1, 1))  
        (4): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (5): ReLU(inplace=True)  
    )  
    (prep1): Sequential(  
        (0): Conv2d(32, 32, kernel_size=(1, 1), stride=(1, 1), groups=32)  
        (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (2): ReLU(inplace=True)  
        (3): Conv2d(32, 32, kernel_size=(1, 1), stride=(1, 1))  
        (4): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (5): ReLU(inplace=True)  
    )  
    (prep2): Sequential(  
        (0): Conv2d(32, 32, kernel_size=(1, 1), stride=(1, 1), groups=32)  
        (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (2): ReLU(inplace=True)  
        (3): Conv2d(32, 16, kernel_size=(1, 1), stride=(1, 1))  
        (4): BatchNorm2d(16, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
        (5): ReLU(inplace=True)  
    )  
Python
```

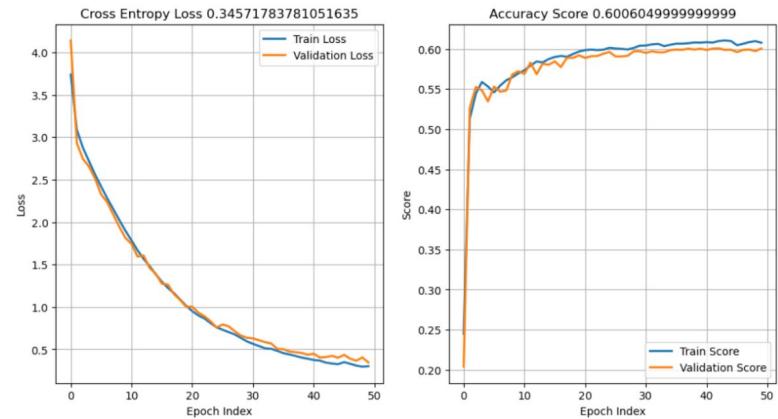


Lets change the network...

UNetSmall CamVid_80_c3

```
trainning   : EPOCHS_TO_RUN=50 ; batch=8
model       : file=./unet_model_large_x80_c3_e50_D.pth
input data  : train_len = 369 ; test_len=100 ; val_len=232
img data    : IMAGE_SIZE=(80, 80)
model       : NUM_CLASS=4 ; NUM_CHANN=3 ; FOLD_RATIO=4
score        : val_score=0.6006049999999999 ; loss_score 0.34571783781051635
Trainable parameters: 41436
parameters_size_in_megabytes : 0.1580657958984375
```

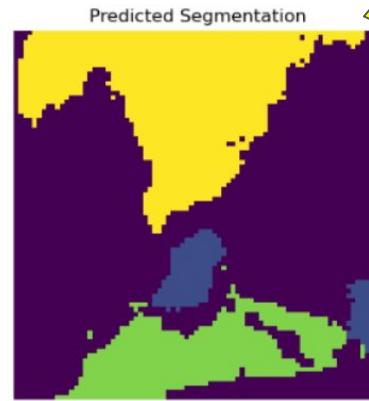
B



Input Image



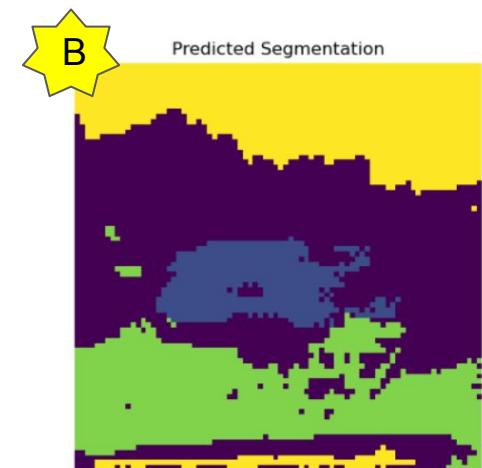
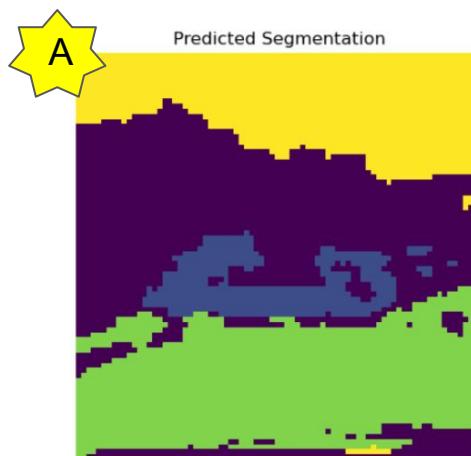
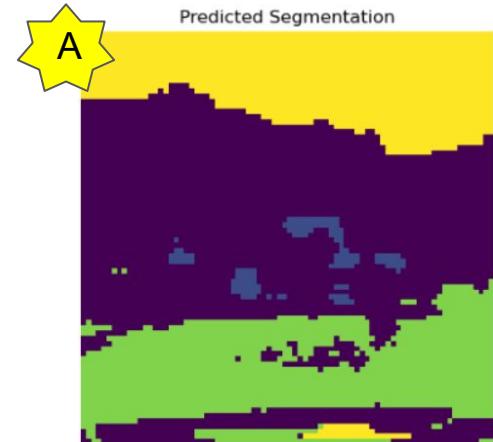
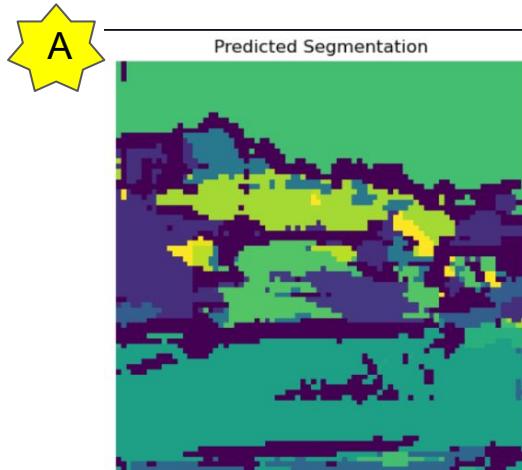
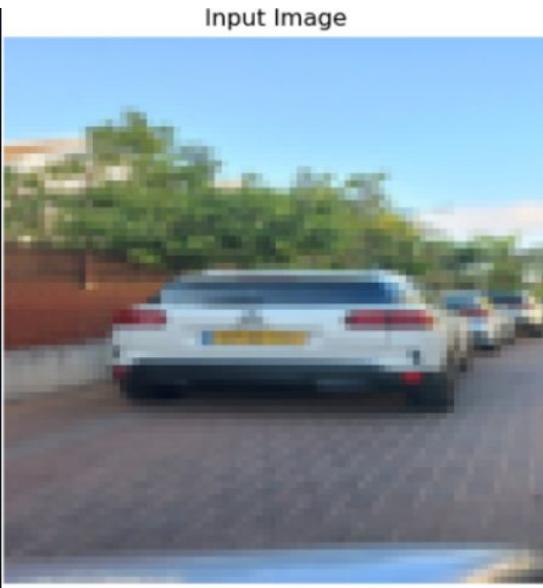
Ground Truth



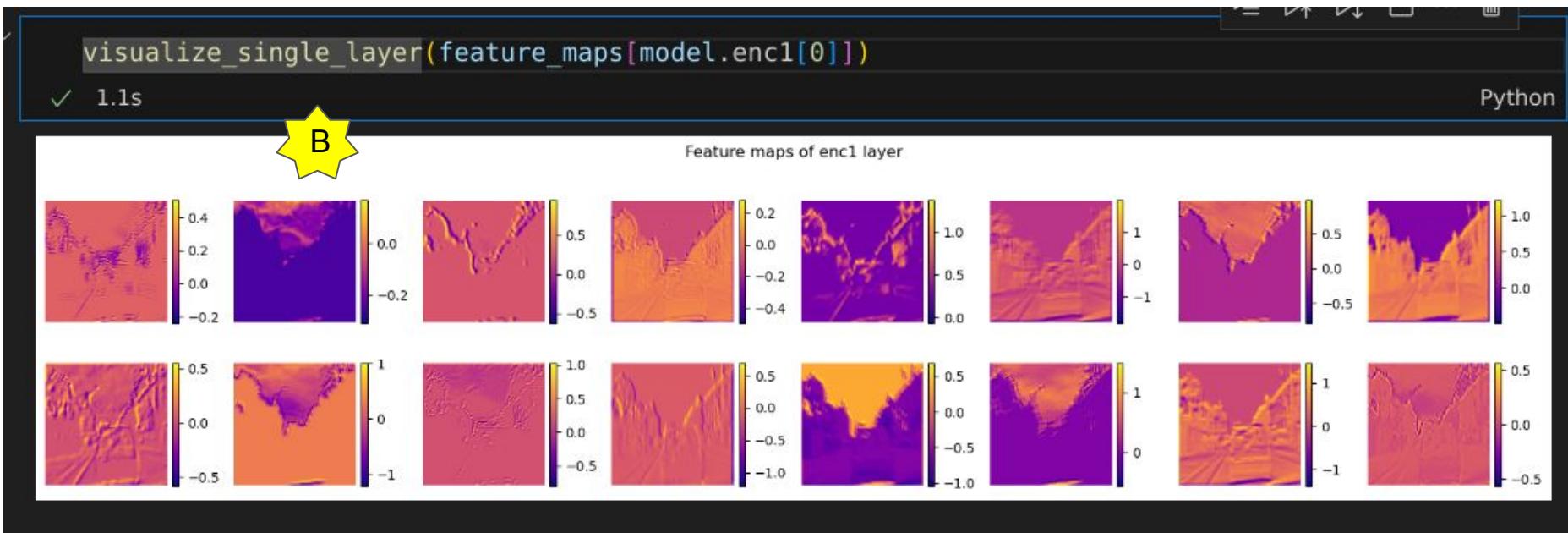
Predicted Segmentation

B

(my) Blind test...



hooks



Make the magic

```
(ai8x-training)python train.py --deterministic --epochs 50  
--optimizer Adam --lr 0.001 --wd 0  
  
--model ai85unetlarge --out-fold-ratio 4 --use-bias  
  
--dataset CamVid_s352_c3 --device MAX78000  
  
--batch-size 16 --qat-policy policies/qat_policy_  
camvid.yaml --compress policies/schedule-camvid.yaml  
--validation-split 0 --compiler-mode none;
```

```
(ai8x-synthesis) python quantize.py  
..../test2_unet/best.pth.tar  
  
..../test2_unet/best-q.pth.tar --device MAX78000 -v
```

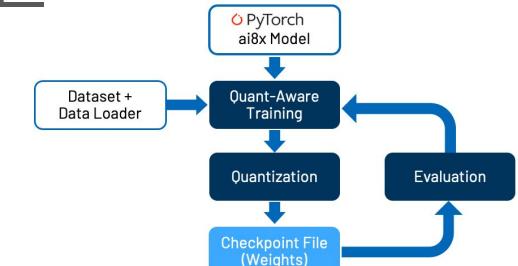
```
(ai8x-training) python train.py --deterministic --model ai85unetlarge --out-fold-ratio 4 --dataset CamVid_s352_c3  
--device MAX78000 --qat-policy policies/qat_policy camvid.yaml --use-bias --evaluate -8 --exp-load-weights-from  
..../test2_unet/best-fakept-q.pth.tar --save-sample 10 --compiler-mode none ;
```

```
(ai8x-synthesize) python ai8xsize.py --verbose --test-dir ..../test2_unet/demo --prefix ai85-unet --checkpoint-file  
..../test2_unet/best-fakept-q.pth.tar --config-file ..../test2_unet/camvid-unet-large-fakept.yaml --device MAX78000  
--overlap-data --mlator --no-unload --max-checklines 8192 --new-kernel-loader --overwrite
```

```
(sdk)>make
```

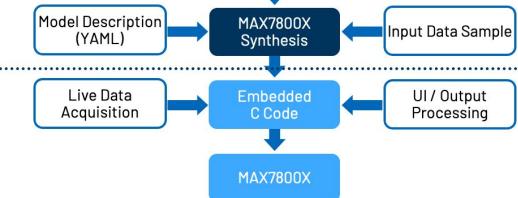
```
(Burn)openocd -s /data/tinyai/maxCNN/sdk/Tools/OpenOCD/scripts -f interface/cmsis-dap.cfg -f target/max78000.cfg -c  
"program build/max78000.elf verify reset exit" ; > minicom
```

① Training ai8x-training

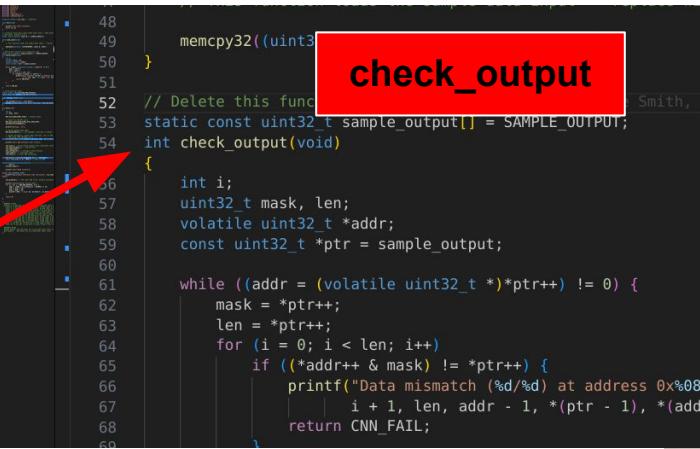


② Synthesis ai8x-synthesis

③ Deployment Embedded SDK

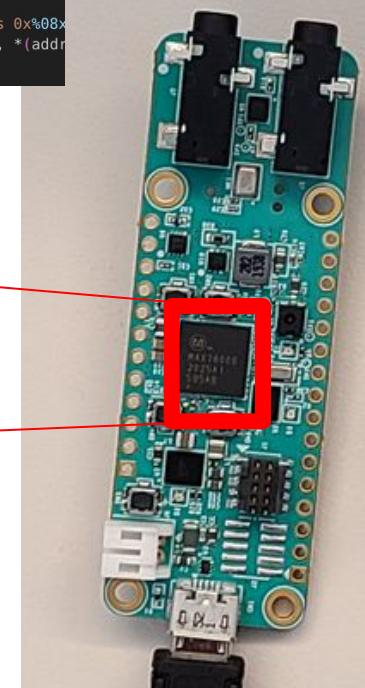
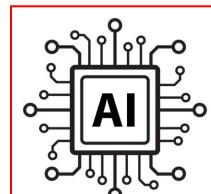
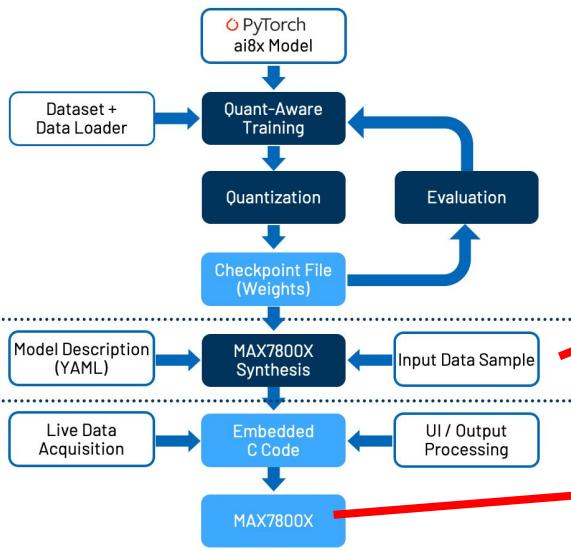


Verify:



```
101 // Enable peripheral, enable CNN interrupt, turn on CNN
102 // CNN clock
103 cnn_enable(); // Set PCLKSEL_PCLK, MXC_S_GCR
104
105 printf("\n***\n");
106
107 cnn_init(); // Bring state machine into consistent state
108 cnn_load_weights(); // Load kernels
109 cnn_load_bias();
110 cnn_configure(); // Configure state machine
111 load_input(); // Load data input
112 cnn_start(); // Start CNN processing
113
114 SCB->SCR &= ~SCB_SCR_SLEEPDEEP_Msk; // SLEEPDEEP=0
115 while (cnn_time == 0) __WFI(); // Wait for CNN
116
117 if (check_output() != CNN_OK)
118     fail();
119 softmax_layer();
120
121 printf("\n*** PASS ***\n\n");
122
```

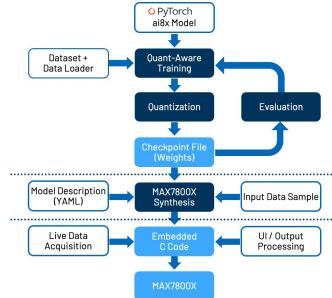
```
48
49     memcpy32(uint32_t *addr, const uint32_t *ptr, uint32_t len);
50 }
51 // Delete this function
52 static const uint32_t sample_output[] = SAMPLE_OUTPUT;
53 int check_output(void)
54 {
55     int i;
56     uint32_t mask, len;
57     volatile uint32_t *addr;
58     const uint32_t *ptr = sample_output;
59
60     while ((addr = (volatile uint32_t *)*ptr++) != 0) {
61         mask = *ptr++;
62         len = *ptr++;
63         for (i = 0; i < len; i++)
64             if ((*addr++ & mask) != *ptr++)
65                 printf("Data mismatch (%d/%d) at address 0x%08x\n"
66                         "    | i + 1, len, addr - 1, *(ptr - 1), *(addr
67                         );
68     }
69 }
```



train

A

```
2024-07-09 06:50:24,177 - ==> Best [Top1: 80.296 Params: 278176 on epoch: 48]
2024-07-09 06:50:24,177 - Saving checkpoint to: logs/2024.07.09-063349/qat_checkpoint.pth.tar
2024-07-09 06:50:24,228 - --- test -----
2024-07-09 06:50:24,228 - 600 samples (16 per mini-batch)
2024-07-09 06:50:26,400 - Test: [ 10/ 38] Loss 0.675823 Top1 79.704277
2024-07-09 06:50:28,607 - Test: [ 20/ 38] Loss 0.673408 Top1 79.728353
2024-07-09 06:50:30,235 - Test: [ 30/ 38] Loss 0.689856 Top1 78.895035
2024-07-09 06:50:31,204 - Test: [ 38/ 38] Loss 0.694540 Top1 78.687414
2024-07-09 06:50:31,331 - ==> Top1: 78.687 Loss: 0.695
```



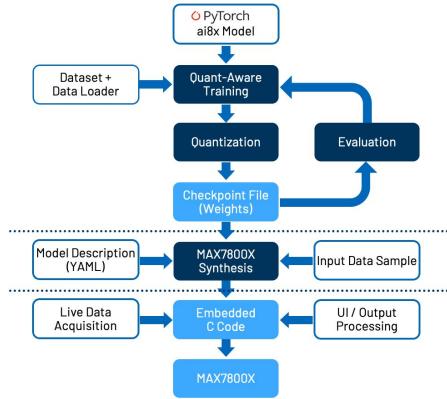
B

```
==> Best [Top1: 77.720 Params: 40504 on epoch: 45]
Saving checkpoint to: logs/2024.07.13-062820/qat_checkpoint.pth.tar
--- test -----
600 samples (16 per mini-batch)
Test: [ 10/ 38] Loss 0.863710 Top1 69.603670
Test: [ 20/ 38] Loss 0.864482 Top1 69.490563
Test: [ 30/ 38] Loss 0.871459 Top1 69.198025
Test: [ 38/ 38] Loss 0.882090 Top1 68.698626
--> Top1: 69.600 Loss: 0.882
```

ensemble

quant

- 32-bit floating-point numbers to 8-bit integers
- Reduces Model Size / Speeds Up Inference / Lowers Power Consumption
While
...retain the model's accuracy as much as possible....



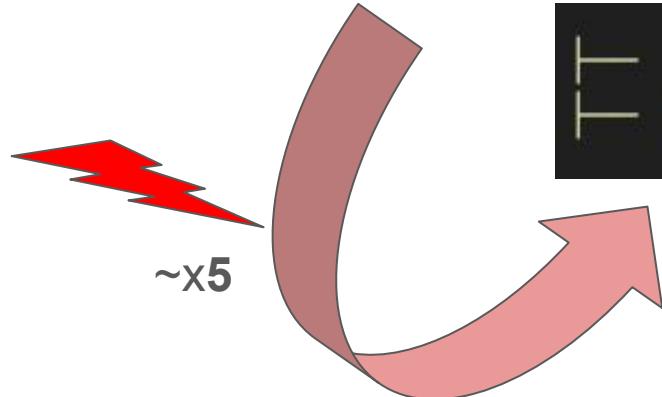
A

3.3M qat_best.pth.tar
3.4M best-q.pth.tar



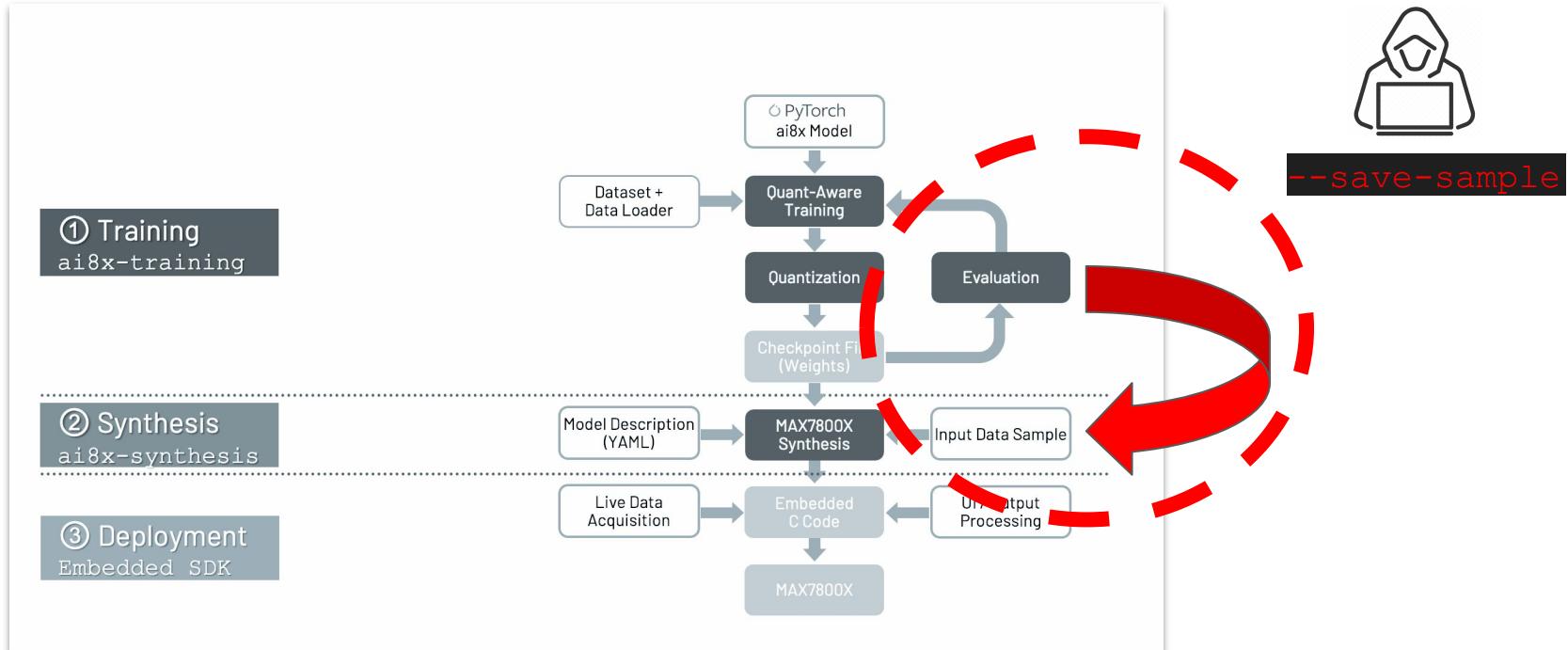
B

[674K] qat_best.pth.tar
[655K] best-q.pth.tar



eval + sample

```
Dataset sizes:  
    test=600  
--- test ---  
600 samples (256 per mini-batch)  
==> Saving sample at index 10 to sample.camvid_s352.c3.npy  
[1] 215756 killed    python train.py --deterministic --model myunetsmall --out-fold-ratio 4  
(ai8x-training) (base) → ai8x-training git:(develop) ✘
```



build

```
#include camra.h
```



(A) 3.4Mb

A

```
|- ai85-unet.launch  
|- cnn.c  
|- cnn.h  
|- log.txt  
└── main.c  
└── Makefile  
└── project.mk  
└── sampledata.h  
└── sampleoutput.h  
└── softmax.c  
└── weights.h
```

RESOURCE USAGE

Weight memory: 281,312 bytes out of 442,368 bytes total (63.6%)
Bias memory: 908 bytes out of 2,048 bytes total (44.3%)

```
if (check_output() != CNN_OK)
| fail();
softmax_layer();

printf("\n*** PASS ***\n\n");
```

$\sim x5$ B (B) 655 K

| max78000.elf

run1.ipynb run5.ipynb x

Select a Kernel from 192.168.2.109 on win

1_unetSmall > run5.ipynb > UNetLarge (7800 to x86) | + Code + Markdown | Run All Clear All Output ★ Python 3 (ipykernel) ~/miniconda3/envs/WorkshopCUDAEnv/bin/python Recommended

UNetLarge (7800 to x86) - CamVid_80_c3

- with feature map
- 80X80 px
- 3 labels

Src code on local host

Remote kernel

```

import torch
import torch.nn as nn
from torch import Module
import torch.optim as optim
from torch.optim import Optimizer

```

PROBLEMS 135 OUTPUT PORTS JUPYTER COMMENTS DEBUG CONSOLE TERMINAL

```

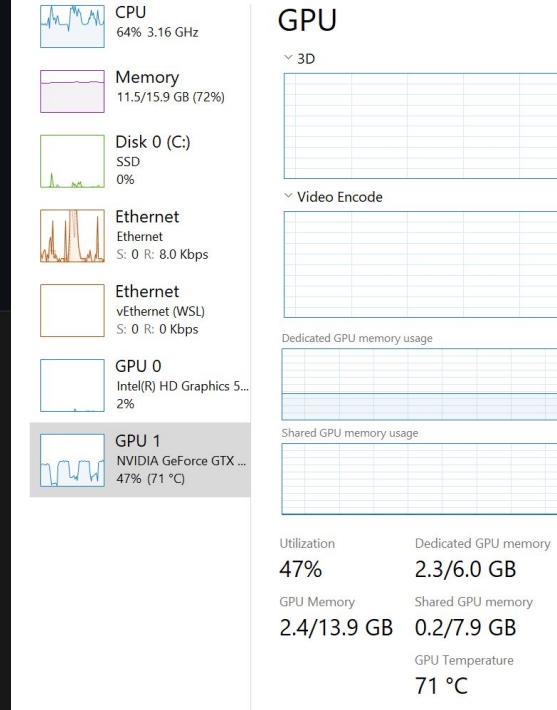
https://ubuntu.com/engage/secure-kubernetes-at-the-edge
Last login: Fri Jul 12 22:49:10 2024 from ::1
#####
###      solov wincpc ubuntu      #####
> conda activate WorkshopCUDAEnv
#####
(base) ~
(base) ~
(base) ~ conda activate WorkshopCUDAEnv
(WorkshopCUDAEnv) ~ jupyter server
[I 2024-07-13 00:27:22.795 ServerApp] Extension package jupyterlab took 0.1220s to import.
[I 2024-07-13 00:27:22.982 ServerApp] jupyter_lsp | extension was successfully linked.
[I 2024-07-13 00:27:22.987 ServerApp] jupyter_server_terminals | extension was successfully linked.
[I 2024-07-13 00:27:22.991 ServerApp] jupyterlab | extension was successfully linked.
[I 2024-07-13 00:27:22.996 ServerApp] notebook | extension was successfully linked.
[I 2024-07-13 00:27:23.385 ServerApp] notebook_shim | extension was successfully linked.
[I 2024-07-13 00:27:23.407 ServerApp] notebook_shim | extension was successfully loaded.

```

ssh

venv

J.Srv

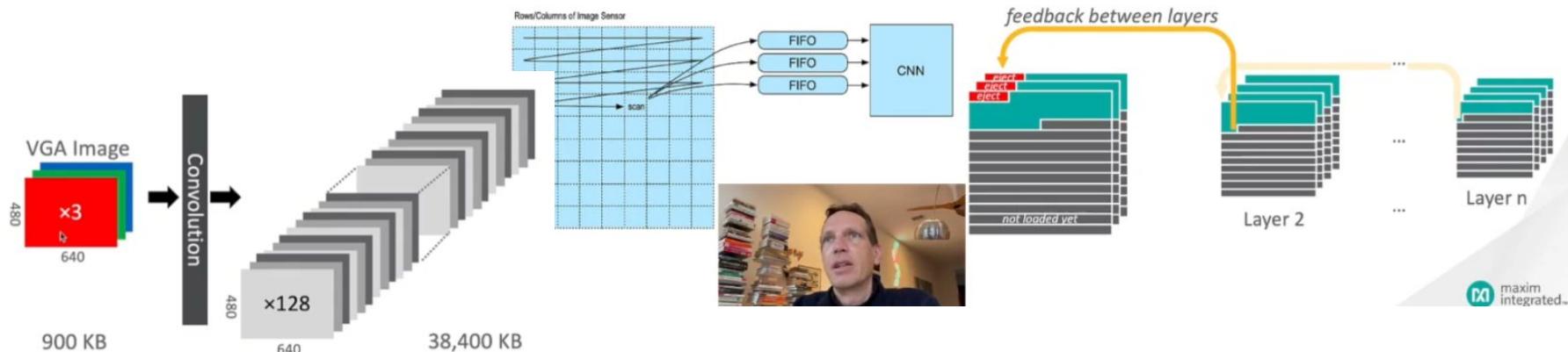


Next?

- Embedded app (camera / sd card / debug port)
 - Measure power / frame rate
 - Larger/multiple networks
 - mobileNet with Unet
-
1. Model Pruning (Structured/Unstructured -close to zero)
 2. Quantization-Aware Training (during training to achieve higher accuracy)
 3. Teacher-Student (smaller network to mimic the outputs of a larger)
 4. Tensor Decomposition
 5. Group Convolutions (limit connections between input and output)
 6. NAS techniques (automatically find a more efficient network architecture)
 7. Weight Sharing (share weights across different layers)
 8. Regularization Techniques (L1/Dropout)
 9. Custom Loss Functions (penalize model size)

Streaming

- Streaming allows input data dimensions that exceed the available per-channel data memory in the accelerator
 - > On MAX78000: anything bigger than 181×181
- Natural when input data is acquired from an image sensor that scans row-by-row



Resolution	Reference	Memory for 128 ch. (max. intermediate)
32x32	CIFAR-10	128 KB
160x160	Intel Movidius	3,200 KB
224x224	ImageNet	6,272 KB
320x240	QVGA	9,600 KB
640x480	VGA	38,400 KB