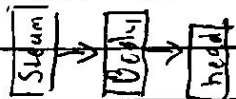# Designing NETWORK Design Spaces

So these guys come in and say hey up until now we have been looking at a model space then finding the best singular model ether Manually or automaticly within that space. But they say lets take it a step/layer above And find a space which describes a subset of all models where these models on avg Better But also "simpler, work well and generalize across settings"

They start with an rather unconstrained space and progressivly constraing it "while maintaing or imporoving" the error distribution produced by the models.

The least constrained space is called AnyNET and is as follows:
(note look at paper for good drawings)

 The Body is where they will Be defining the model + the Bulk of the work will Be done. Stem = Stide2 3×3 conv the 32 channels + the head is AvG Pooling + a FC layer. There are 4 stages in the Body Where each stage i has $b_i$ (blocks) $w_i$ (width) + other Block params.

Since each network has 4 stages
+ each stage has 4 degrees of
freedom in total there are 16 degrees
of freedom.

each stage has $d_i$ (Blocks)
$w_i$ (width), $b_i$ (Bottleneck ratio), and $g_i$ (group width)

$$d_i \leq 16 \qquad w_i \leq 1024 \text{ (and divisable by 8)}$$

$$b_i = \{1, 2, 4\} \qquad g_i \in \{1, 2, \ldots 32\}$$

So above is the AnyNet design space
with $10^{18}$ possible model configs.

Step one they set all $b_i = b$
so it's the same across all stages of
a model. They find no increas in
error but now the design space
is simpler.

Step two is to set $g_i = g$
as above + they find same result

Step 3 they find pattern after
Step 2 where increasing width over the
stages results in Better models
so they test AnyNet$_C$ where AnyNet$_C$
is after step 2 + AnyNet$_D$ is
after step 2 + only models where
$w_{i+1} \geq w_i$ + find it significantly
Better distribution of error

Step 4 they find that similar
as with step 3 if now we
increase depth $d_{c+1} \geq d_c$ the models
are better.

So after all these reductions
our design space went from
$10^{18}$ posobility to $10^7$

So then they come up with the
final design space described as such:

RegNET generated from: $d$, $W_0$, $W_a$, $W_m$
$d < 64$     $W_0, W_a \leq 256$ But we have
                        $\swarrow$ to dicintize    Control the
                                                    scaling of
$1.5 \leq W_m \leq 3$    $U_j = W_0 + W_a \cdot j$ via $W_m$    width

The original tests we have Been
reading about are all done in the
Low epochs + low compute range.
So now they compare in higher compute
higher epoch + 5 stages
the ordering is always RegNet $\nearrow$ AnyNet$_E$ $\nearrow$ AnyNet$_A$
                                                    $\nwarrow$        $\uparrow$
                                                    Post        Pre
                                                    step4      step1

They then have further observations
that the common $b < 1$ + $g = 1$
are not as good as $b = 1$ + $g \geq 1$

Al.

they also found optimal depth = 20
Blocks (interesting Deeper not always Better!)

+ a width multiple of 2.5 (close to the
common one of 2

So Now lets compare ResNET
model to other models.

RegNET Models tend to have
lower Flops But maintain or
Better results the ResNET

In general the RegNETS
Matched or did Better than
state of the art Res NET

and at low flop Efficent net
Better But at Higher flops RegNET
Better. + is much faster in the
Higher flop regions

$d = 4$

Say: $W_0 = 32$

$W_a = 8$

$W_m = 2$

ON quantization: So we get Powers of 2

| | | | |
|---|---|---|---|
| $32 = U_0 = 32 + 8 \cdot 0$ | $\log\left(\frac{32}{32}\right) = 0$ | Round $= 0$ | $W_0 = 32 \cdot 2^0 = 32$ |
| $40 = U_1 = 32 + 8 \cdot 1$ | $\log(40/32) = .32$ | $\Rightarrow = 0$ | $W_1 = 32 \cdot 2^0 = 32$ |
| $48 = U_2 = 32 + 8 \cdot 2$ | $\log(48/32) = .58$ | $= 1$ | $W_2 = 32 \cdot 2^1 = 64$ |
| $56 = U_3 = 32 + 8 \cdot 3$ | $\log(56/32) = .81$ | $= 1$ | $W_3 = 32 \cdot 2^1 = 64$ |

Now 2 Stages

| | Stage 1 | Stage 2 |
|---|---|---|
| | 2 Blocks | 2 Blocks |
| | 32 width | 64 width |