# QANet.

## Model



start prob.  →  end prob.

softmax  ←  softmax

linear  ←  linear

concat  ←  concat

Stacked Model Encoder Blocks

SMEB

SMEB

Context - Query Attention

E | SMEB    E | SMEB    E: embedding

Embedding    Embedding

[O O O O]    [O O O O]

context      Question.

**Model Layers:**

1. Input Embedding L.
2. Embedding Encoder L.
3. Content - Query Attention L.
4. Model Encoder L.
5. Output L.

---

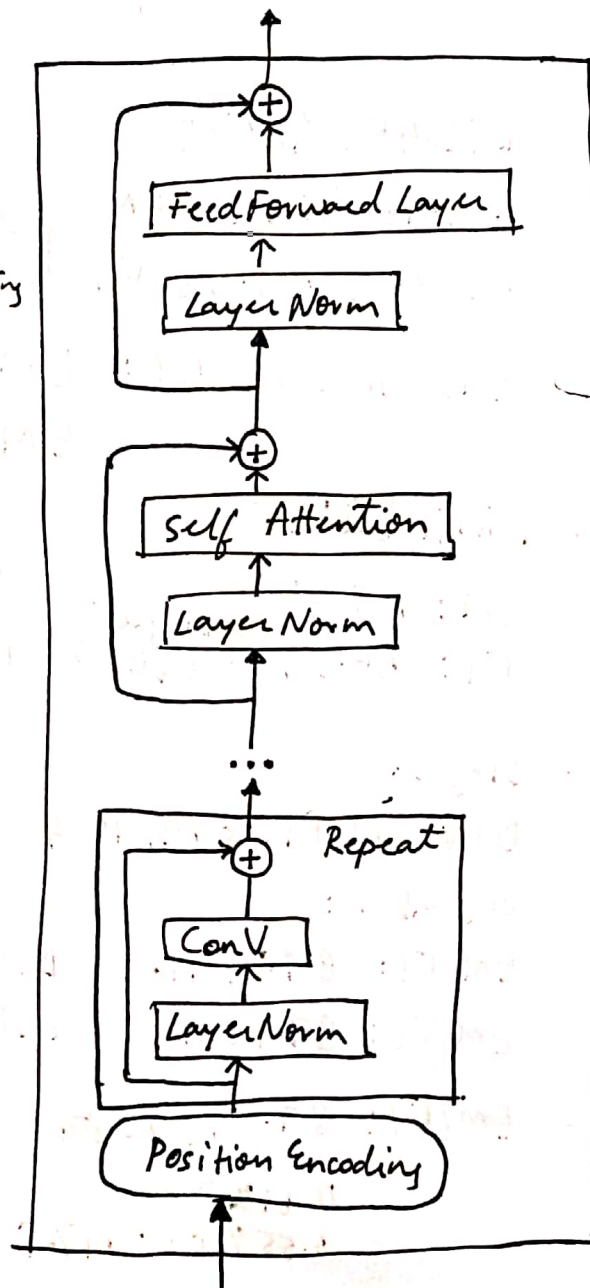- 77.0 F1 score → 3 hrs training (BiDAF: 15 hrs)
- 82.7 F1 (dev) → 18 hrs.

---

- SQuAD → 84.6 F1 (test) (81.8 SoTA 2017)

- 2018:
  SQuAD     EM/F1
  82.2/88.6 (single)
  83.9/89.7 (ensemble)
  82.3 (human)  ●

---

## One Encoder Block.



⊕

Feed Forward Layer

Layer Norm

⊕

Self Attention

Layer Norm

⋯

⊕  Repeat

Conv

Layer Norm

Position Encoding

Scanned with CamScanner

# QANET

context paragraph : $C = \{c_1, ..., c_n\}$

query sentence : $Q = \{q_1, ..., q_m\}$

o/p span : $S = \{c_i, ..., c_{i+j}\}$

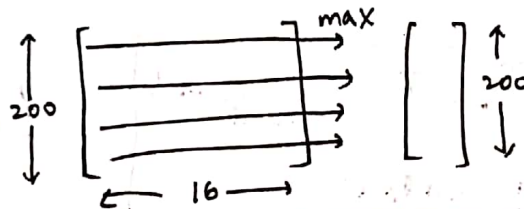$\boxed{x \in C, Q}$

① $\boxed{\text{Input Embedding Layer :}}$

- **word emb**: fixed
  - $P_1 = 300$ D (GloVe)
  - $\langle unk \rangle \longrightarrow$ trainable $\bar{c}$ random init

- char emb :
  - $P_2 = 200$ D
  - max-word-len = 16
  - $\langle pad \rangle$

$\boxed{x = [x_w; x_c] \in R^{P_1 + P_2} \in R^{500}}$

② $\boxed{\text{Embedding Encoder Layer}}$

- stack of the foll$^n$ block :
  $[\text{conv layer} \times \# + \text{self attention layer} + \text{feed forward layer}]$

- depthwise separable conv. (mem. efficient)
  - $k = 7$
  - $d = 128$ filters
  - \# conv layers within 1 block = 4.  $\Big\}$ conv.

- multiheaded attention
  - query, key      ~~inside a residual block.~~
  - no. of heads = 8  $\Big\}$ self atten.

- conv / attent / feed forw
  - inside a residual block
  - $f(\text{layer norm}(x) + x)$
  - no of encoder blocks = 1  $\Big\}$ General

- i/p = 500 D $\longrightarrow$ 1D conv $\longrightarrow$ 128 D *
  o/p = 128 D  $\Big\}$ i/p o/p

③ **Context - Query Attention Layer.**

- Similarity matrix $S \in R^{n \times m}$ : sim b/w $C$ & $Q$.
  
  ($n$) ($m$)

- $\boxed{S \xrightarrow{softmax} \bar{S}}$ (row normalized)

  $\Rightarrow dim = 1$

- Attention $(C \to Q)$:

  $\boxed{A = \bar{S} \cdot Q^T \in R^{n \times d}}$
  
  $(n \times m)$ $(m \times d)$

- Similarity fn : (trilinear)

  $\boxed{f(q, c) = W_0 [q, c, q \odot c]}$

---

✳ **Query - Context Attention.**

- (DCN)

- column normalized $\boxed{S \xrightarrow{soft} \bar{\bar{S}}}$ $\Rightarrow dim = 0$

  $\boxed{B = \bar{S} \cdot \bar{\bar{S}}^T \cdot C^T \in R^{n \times d}}$   ? $(m \times d)$

  $(n \times m)$ $(m \times n)$ $(n \exists \times d)$

④ **Model Encoder Layer**

- i/p : $[c, a, c \odot a, c \odot b]$     $a, b$: row of $A$ & $B$

- parameters same as before except:

  #conv layers = 2  in a block

  # blocks = 7

  $\Rightarrow$ weights are shared b/w each of the 3 repitions of this Layer.

⑤ **Output Layer**

$\boxed{p^1 = softmax(W_1 [M_0; M_1])  \quad p^2 = softmax(W_2 [M_0; M_2]).}$

$M_2$  
$M_1$  
$M_0$

$\boxed{score = p_s^1 \times p_e^2}$   $\boxed{L(\theta) = -\frac{1}{N} \sum_i^N [log(p^1 y_i^1) + log(p^2 y_i^2)]}$

Inference : $p_s^1 p_e^2 \longrightarrow max$   $s \leq e$.

(DP can solve in linear time)

# EXPERIMENTS

(SQuAD 1)

- **Data Preprocessing**
  - NLTK tokenizer
  - max context len = 400
  - (para longer → discard)
  - max ans len = 30
  - (300-D GLoVe)

- $<PAD>$ : short
- $<UNK>$ → trained
- 200 D → trained char emb

- **Training Details**
  - ① L2 weight decay
    - → all trainable variables
    - → $\lambda = 3 \times 10^{-7}$
  - ② Dropout
    - → word emb ⟩ 0.1
    - → char emb → 0.05
    - → blw layers → 0.1

- hidden size ⟩ 128
- conv filter ⟩

- batch size = 32   steps = 150 K
- conv : emb = 4   k = 7   # = 1
- conv : model = 2   k = 5   # = 7

- **Stochastic Depth Method** (layer dropout)
  - → within each embedding / model encoder layer.
  - → sublayer $l$ : survival prob

  $$P_l = 1 - \frac{l}{L}(1 - P_L)$$

  $P_L = 0.9$    $L$ = last layer.

- **ADAM**
  - $\beta_1 = 0.8$    $\varepsilon = 10^{-7}$
  - $\beta_2 = 0.999$
  - LR warm up scheme
    - $0.0 \rightarrow 0.001$   (1000 steps)
    - then constant.

- Exponential moving avg is applied on all trainable variables decay rate = 0.9999.

NVIDIA P100 GPU.

no data aug :   EM / FI   73.6 / 82.7

(SQuAD 1)