# ① Neural Turing Machine

**Reading** : memory $M_t \to$ $\underset{\text{memory location}}{\boxed{N}} \times \underset{\text{vector size at each location}}{\boxed{M}}$ matrix

(with $\downarrow$ time)

Attention weight, $0 \leq w_t(i) \leq 1$ ; $i \sim \{1 \dots N\}$ with constraint $\boxed{\sum_i w_t(i) = 1}$

read memory $\underset{(1 \times m)\text{size}}{} \quad r_t \leftarrow \boxed{\sum_i w_t(i) \, M_t(i)} \to$ i'th row

$\to$ convex combination. weysum of rows,

Differentiable ,

**Writing** :

$[1 \, 1 \cdots 1] \to$ erase vector $(1 \times m)$

$$\tilde{M}_t(i) \leftarrow M_{t-1}(i) \underbrace{\left[ 1 - w_t(i) \, e_t \right]}_{\text{pointwise multiply}}$$

$\downarrow$ After erase add

$\to$ Both Differentiable

$$M_t(i) \leftarrow \tilde{M}_t(i) + w_t(i) \, a_t$$

? Constructing weight vector!

focusing by content

key strength.
cosine sim
$\uparrow$ key vector $(1 \times m)$

$$w_t^c(i) \leftarrow \frac{\exp\left( \beta_t k \left[ k_t , M_t(i) \right] \right)}{\sum_j \exp\left( \beta_t k \left[ k_t , M_t(j) \right] \right)}$$

focusing by location:

interpolation gate $(0, 1)$

$$\underline{w}_t^g \leftarrow g_t \underline{w}_t^c + (1 - g_t) \underline{w}_{t-1}$$

$$\tilde{w}_t(i) \leftarrow \sum_{j=0}^{N-1} w_t^g(j) \, s(i-j)$$

$\underbrace{\qquad}$ shift weight.

$$w_t(i) \leftarrow \frac{\tilde{w}_t(i)^{\gamma_t}}{\sum_i \tilde{w}_t(i)^{\gamma_t}} \Bigg] \rightarrow \text{sharpening} \, \& \, \text{Normalization.}$$