

used to model fine-grained time scales.

### 10.9.1 Adding Skip Connections through Time

One way to obtain coarse time scales is to add direct connections from variables in the distant past to variables in the present. The idea of using such skip connections dates back to [Lin et al. \(1996\)](#) and follows from the idea of incorporating delays in feedforward neural networks ([Lang and Hinton, 1988](#)). In an ordinary recurrent network, a recurrent connection goes from a unit at time  $t$  to a unit at time  $t + 1$ . It is possible to construct recurrent networks with longer delays ([Bengio, 1991](#)).

As we have seen in section 8.2.5, gradients may vanish or explode exponentially *with respect to the number of time steps*. [Lin et al. \(1996\)](#) introduced recurrent connections with a time-delay of  $d$  to mitigate this problem. Gradients now diminish exponentially as a function of  $\frac{\tau}{d}$  rather than  $\tau$ . Since there are both delayed and single step connections, gradients may still explode exponentially in  $\tau$ . This allows the learning algorithm to capture longer dependencies although not all long-term dependencies may be represented well in this way.

### 10.9.2 Leaky Units and a Spectrum of Different Time Scales

Another way to obtain paths on which the product of derivatives is close to one is to have units with *linear* self-connections and a weight near one on these connections.

When we accumulate a running average  $\mu^{(t)}$  of some value  $v^{(t)}$  by applying the update  $\mu^{(t)} \leftarrow \alpha\mu^{(t-1)} + (1 - \alpha)v^{(t)}$  the  $\alpha$  parameter is an example of a linear self-connection from  $\mu^{(t-1)}$  to  $\mu^{(t)}$ . When  $\alpha$  is near one, the running average remembers information about the past for a long time, and when  $\alpha$  is near zero, information about the past is rapidly discarded. Hidden units with linear self-connections can behave similarly to such running averages. Such hidden units are called **leaky units**.

Skip connections through  $d$  time steps are a way of ensuring that a unit can always learn to be influenced by a value from  $d$  time steps earlier. The use of a linear self-connection with a weight near one is a different way of ensuring that the unit can access values from the past. The linear self-connection approach allows this effect to be adapted more smoothly and flexibly by adjusting the real-valued  $\alpha$  rather than by adjusting the integer-valued skip length.

These ideas were proposed by [Mozer \(1992\)](#) and by [El Hihhi and Bengio \(1996\)](#). Leaky units were also found to be useful in the context of echo state networks ([Jaeger et al., 2007](#)).