

Change in the value function due to learning about a particular  $b'$ :

$$\begin{aligned} v(b') - v(b) &= \sum_a \pi(a | b') q(b', a) - \sum_{a'} \pi(a | b) q(b, a) \\ &= \sum_a [(\pi(a | b') - \pi(a | b)) q(b', a) + \pi(a | b) (q(b', a) - q(b, a))] \end{aligned} \quad (1)$$

$$\begin{aligned} q(b', a) - q(b, a) &= \sum_{b''} p(b'' | b', a) [r(b', a) + \gamma v(b'')] \\ &\quad - \sum_{b'} p(b' | b, a) [r(b, a) + \gamma v(b')] \\ &= r(b', a) + \gamma \sum_{b''} p(b'' | b', a) v(b'') \\ &\quad - r(b, a) + \gamma \sum_{b'} p(b' | b, a) v(b') \\ &= \underbrace{r(b', a) - r(b, a)}_{\text{Difference in the expected immediate return}} + \underbrace{\gamma \left[ \sum_{b''} p(b'' | b', a) v(b'') - \sum_{b'} p(b' | b, a) v(b') \right]}_{\text{Difference in the expected future return}} \end{aligned} \quad (2)$$

Note that we can write

$$\begin{aligned} cv(b'') - dv(b') &= cv(b'') - cv(b') + cv(b') - dv(b') \\ &= c(v(b'') - v(b')) + v(b')(c - d) \end{aligned} \quad (3)$$

Therefore

$$\begin{aligned} q(b', a) - q(b, a) &= r(b', a) - r(b, a) + \gamma \left[ \sum_{b''} p(b'' | b', a) v(b'') - \sum_{b'} p(b' | b, a) v(b') \right] \\ &= r(b', a) - r(b, a) + \gamma \sum_{b''} p(b'' | b', a) [v(b'') - v(b')] \\ &\quad + \gamma v(b') \left[ \sum_{b''} p(b'' | b', a) - \sum_{b'} p(b' | b, a) \right] \end{aligned} \quad (4)$$

Note that the last term goes to zero. Therefore, substituting in:

$$\begin{aligned} v(b') - v(b) &= \sum_a \pi(a | b') q(b', a) - \sum_{a'} \pi(a | b) q(b, a) \\ &= \sum_a [\pi(a | b') - \pi(a | b)] q(b', a) \\ &\quad + \sum_a \pi(a | b) [r(b', a) - r(b, a)] \\ &\quad + \gamma \sum_a \pi(a | b) \left[ \sum_{b''} p(b'' | b', a) (v(b'') - v(b')) \right] \end{aligned}$$

Unrolling:

$$\begin{aligned}
v(b') - v(b) = & \sum_{i=0}^{\infty} \sum_{b' \in \mathcal{B}} \gamma^i P(b \rightarrow b', i, \pi(b)) \times \\
& \sum_a \left( \underbrace{[\pi(a | b') - \pi(a | b)] q(b', a)}_{\text{Localised Gain}} + \underbrace{\mathbb{E}_{\pi(b)} [r(b', a) - r(b, a)]}_{\text{Accumulates into long-term consequences}} \right)
\end{aligned} \tag{5}$$

Equation 5 thus shows the non-local effect of policy change at a single belief state. Even though Gain is only non-zero at the exact (belief) location where we perform an update, this update nonetheless has long-lasting (discounted) consequences.