Change in the value function due to learning after taking action $a^*$:

$$v(ba^*) - v(b) = \sum_{b'} p(b' \mid b, a^*) \sum_a \pi(a \mid b') q(b', a) - \sum_a \pi(a \mid b) q(b, a) \tag{1}$$

$$= \sum_a \left( \sum_{b'} p(b' \mid b, a^*) \pi(a \mid b') q(b', a) - \pi(a \mid b) q(b, a) \right)$$

$$= \sum_a \left( \sum_{b'} p(b' \mid b, a^*) \pi(a \mid b') q(b', a) - \pi(a \mid b) q(b, a) \right.$$

$$\left. + \pi(a \mid b) \sum_{b'} p(b' \mid b, a^*) q(b', a) - \pi(a \mid b) \sum_{b'} p(b' \mid b, a^*) q(b', a) \right)$$

$$= \sum_a \left( \sum_{b'} p(b' \mid b, a^*) \big( \pi(a \mid b') - \pi(a \mid b) \big) q(b', a) \right.$$

$$\left. + \pi(a \mid b) \sum_{b'} p(b' \mid b, a^*) \big( q(b', a) - q(b, a) \big) \right)$$

$$= \sum_{b'} p(b' \mid b, a^*) \sum_a \left( \big( \pi(a \mid b') - \pi(a \mid b) \big) q(b', a) \right.$$

$$\left. + \pi(a \mid b) \big( q(b', a) - q(b, a) \big) \right)$$

$$\tag{2}$$

$$q(b', a) - q(b, a) = \sum_{b''} p(b'' \mid b', a) \big[ r(b', a) + \gamma v(b'') \big] \tag{3}$$

$$- \sum_{b'} p(g' \mid b, a) \big[ r(b, a) + \gamma v(b') \big]$$

$$= r(b', a) + \gamma \sum_{b''} p(b'' \mid b', a) v(b'')$$

$$- r(b, a) + \gamma \sum_{b'} p(b' \mid b, a) v(b')$$

$$= \underbrace{r(b', a) - r(b, a)}_{\substack{\text{Difference in the expected} \\ \text{immediate return}}} + \underbrace{\gamma \Big[ \sum_{b''} p(b'' \mid b', a) v(b'') - \sum_{b'} p(b' \mid b, a) v(b') \Big]}_{\substack{\text{Difference in the expected} \\ \text{future return}}}$$

1