

**Electronic supplement.** Calculation of some mathematical results for "Reciprocal mimicry": kin selection can drive defended prey to resemble their Batesian mimics".

### A1. Calculation of $\left[ \frac{\partial t_{mix}^*}{\partial m'_d} \right]_{m'_d=\hat{m}_d}$ .

The optimal attack threshold  $t_{mix}^*(m'_d, \hat{m}_d, \hat{m}_m)$  is the unique root of

$$f(t) = rZ'(t - m'_d) + (1 - r)Z'(t - \hat{m}_d) - KZ'(t - \hat{m}_m). \text{ Let } t_{mix}^* \text{ be shorthand for } t_{mix}^*(m'_d, \hat{m}_d, \hat{m}_m).$$

According to the implicit function theorem, we have

$$\frac{\partial t_{mix}^*}{\partial m'_d} = \left[ -\frac{\partial f / \partial m'_d}{\partial f / \partial t} \right]_{t=t_{mix}^*} = -\frac{-rZ''(t_{mix}^* - m'_d)}{rZ''(t_{mix}^* - m'_d) + (1 - r)Z''(t_{mix}^* - \hat{m}_d) - KZ''(t_{mix}^* - \hat{m}_m)},$$

which using the fact that  $Z''(x) = -Z'(x)x$  can be written as

$$\frac{\partial t_{mix}^*}{\partial m'_d} = -\frac{rZ'(t_{mix}^* - m'_d)(t_{mix}^* - m'_d)}{rZ'(t_{mix}^* - m'_d)(-1)(t_{mix}^* - m'_d) + (1 - r)Z'(t_{mix}^* - \hat{m}_d)(-1)(t_{mix}^* - \hat{m}_d) - KZ'(t_{mix}^* - \hat{m}_m)(-1)(t_{mix}^* - \hat{m}_m)}.$$

As by definition  $f(t_{mix}^*) = 0$ , we may substitute in  $rZ'(t_{mix}^* - m'_d) + (1 - r)Z'(t_{mix}^* - \hat{m}_d)$  for

$KZ'(t_{mix}^* - \hat{m}_m)$  and rearrange, yielding

$$\begin{aligned} \frac{\partial t_{mix}^*}{\partial m'_d} &= \frac{rZ'(t_{mix}^* - m'_d)(t_{mix}^* - m'_d)}{rZ'(t_{mix}^* - m'_d)(t_{mix}^* - m'_d - t_{mix}^* + \hat{m}_m) + (1 - r)Z'(t_{mix}^* - \hat{m}_d)(t_{mix}^* - \hat{m}_d - t_{mix}^* + \hat{m}_m)} \\ &= \frac{rZ'(t_{mix}^* - m'_d)(t_{mix}^* - m'_d)}{rZ'(t_{mix}^* - m'_d)(\hat{m}_m - m'_d) + (1 - r)Z'(t_{mix}^* - \hat{m}_d)(\hat{m}_m - \hat{m}_d)}. \end{aligned}$$

We are interested in the case in which  $m'_d = \hat{m}_d$ , in which case  $t_{mix}^*$  is given explicitly by (2.2) in the main text. We obtain

$$\begin{aligned} \left[ \frac{\partial t_{mix}^*}{\partial m'_d} \right]_{m'_d=\hat{m}_d} &= \frac{rZ'(t^* - \hat{m}_d)(t^* - \hat{m}_d)}{rZ'(t^* - \hat{m}_d)(\hat{m}_m - \hat{m}_d) + (1 - r)Z'(t^* - \hat{m}_d)(\hat{m}_m - \hat{m}_d)} \\ &= \frac{r(t^* - \hat{m}_d)}{r(\hat{m}_m - \hat{m}_d) + (1 - r)(\hat{m}_m - \hat{m}_d)} = \frac{r \left( \frac{\hat{m}_m + \hat{m}_d}{2} - \frac{\ln K}{\hat{m}_m - \hat{m}_d} - \hat{m}_d \right)}{d} = r \left( \frac{1}{2} - \frac{\ln K}{d^2} \right). \end{aligned}$$

## A2. Determining the sign of the fitness gradient using attack function $A_{alt}$ .

The rate of attack is given by:

$$A_{alt}(m, m'_d, \hat{m}_d, \hat{m}_m) = \left(1 - Z(t_{mix}^*(m'_d, \hat{m}_d, \hat{m}_m) - m)\right) \left(w_{opt}(m'_d, \hat{m}_d, \hat{m}_m)\right).$$

The two scenarios are treated one after the other.

*Scenario 1.*

The fitness gradient has the same sign as  $-E_r \left[ \frac{\partial}{\partial m'_d} A_{alt}(m'_d, m'_d, \hat{m}_d, \hat{m}_m) \right]_{m'_d=\hat{m}_d}$ .

Multiplying this with the positive expression  $1/(1-p)c$ , we obtain the sign-equivalent expression

$$-E_r \left[ \frac{\partial}{\partial m'_d} \left( 1 - Z(t_{mix}^* - m'_d) \right) \left( (1 - Z(t_{mix}^* - \hat{m}_m))K + rZ(t_{mix}^* - m'_d) + (1-r)Z(t_{mix}^* - \hat{m}_d) - 1 \right) \right]_{m'_d=\hat{m}_d}.$$

Performing the differentiation, we obtain:

$$\begin{aligned} & -E_r \left[ -Z'(t_{mix}^* - m'_d) \left( \frac{\partial t_{mix}^*}{\partial m'_d} - 1 \right) \left( (1 - Z(t_{mix}^* - \hat{m}_m))K + rZ(t_{mix}^* - m'_d) + (1-r)Z(t_{mix}^* - \hat{m}_d) - 1 \right) \right. \\ & \left. + (1 - Z(t_{mix}^* - m'_d)) \left( -Z'(t_{mix}^* - \hat{m}_m)K \frac{\partial t_{mix}^*}{\partial m'_d} + rZ'(t_{mix}^* - m'_d) \left( \frac{\partial t_{mix}^*}{\partial m'_d} - 1 \right) + (1-r)Z'(t_{mix}^* - \hat{m}_d) \frac{\partial t_{mix}^*}{\partial m'_d} \right) \right]_{m'_d=\hat{m}_d}. \end{aligned}$$

When  $m'_d = \hat{m}_d$ , we can (in turn) make the substitutions  $t_{mix}^*(\hat{m}_d, \hat{m}_d, \hat{m}_m) = t^*(\hat{m}_d, \hat{m}_m)$  (see main text),  $Z'(t^* - \hat{m}_m) = Z'(t^* - \hat{m}_d)/K$  (the latter can be checked by substituting in

$$t^* = \frac{\hat{m}_m + \hat{m}_d}{2} - \frac{\ln K}{\hat{m}_m - \hat{m}_d} \text{ and rearranging), and } \left. \frac{\partial t_{mix}^*}{\partial m'_d} \right|_{m'_d=\hat{m}_d} = r \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) \text{ (see section A1).}$$

Making these substitutions and simplifying we obtain

$$\begin{aligned} & \mathbb{E}_r \left[ Z'(t^* - m'_d) \left( \left( \frac{\partial t^*}{\partial m'_d} - 1 \right) ((1 - Z(t^* - \hat{m}_m))K + Z(t^* - \hat{m}_d) - 1) + r(1 - Z(t^* - \hat{m}_d)) \right) \right]_{m'_d = \hat{m}_d} \\ &= Z'(t^* - \hat{m}_d) \mathbb{E}_r \left[ \left( r \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) - 1 \right) ((1 - Z(t^* - \hat{m}_m))K + Z(t^* - \hat{m}_d) - 1) + r(1 - Z(t^* - \hat{m}_d)) \right]. \end{aligned}$$

The fitness gradient has the same sign as the expectation, and by making substitutions

$t^* - \hat{m}_d = d / 2 - \ln(K) / d$  and  $t^* - \hat{m}_m = -d / 2 - \ln(K) / d$  (see main text) and rearranging we obtain (3.2) in the main text:

$$\left( \bar{r} \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) - 1 \right) \left( 1 - Z \left( -\frac{d}{2} - \frac{\ln K}{d} \right) \right) K - \left( \bar{r} \left( -\frac{1}{2} - \frac{\ln K}{d^2} \right) - 1 \right) \left( 1 - Z \left( \frac{d}{2} - \frac{\ln K}{d} \right) \right).$$

*Scenario 2.*

The fitness gradient has the same sign as

$$-\mathbb{E}_{\rho,f} \left[ \rho \frac{\partial}{\partial m'_d} A_{alt}(m'_d, \hat{m}_d, \hat{m}_d, \hat{m}_m) + \frac{\partial}{\partial m'_d} A_{alt}(\hat{m}_d, m'_d, \hat{m}_d, \hat{m}_m) \right]_{m'_d = \hat{m}_d}.$$

Multiplying this with the positive expression  $1 / ((1-p)c)$  and using

$t_{mix}^*(\hat{m}_d, \hat{m}_d, \hat{m}_m) = t^*(\hat{m}_d, \hat{m}_m)$ , we obtain the sign-equivalent expression

$$\begin{aligned} & -\mathbb{E}_{\rho,f} \left[ \rho \frac{\partial}{\partial m'_d} \left( 1 - Z(t^* - m'_d) \right) \left( (1 - Z(t^* - \hat{m}_m))K + rZ(t^* - \hat{m}_d) + (1-r)Z(t^* - \hat{m}_d) - 1 \right) \right. \\ & \left. + \frac{\partial}{\partial m'_d} \left( 1 - Z(t_{mix}^* - \hat{m}_d) \right) \left( (1 - Z(t_{mix}^* - \hat{m}_m))K + rZ(t_{mix}^* - m'_d) + (1-r)Z(t_{mix}^* - \hat{m}_d) - 1 \right) \right]_{m'_d = \hat{m}_d}. \end{aligned}$$

Performing the differentiation, we obtain:

$$\begin{aligned} & -\mathbb{E}_{\rho,f} \left[ \rho Z'(t^* - m'_d) \left( (1 - Z(t^* - \hat{m}_m))K + rZ(t^* - \hat{m}_d) + (1-r)Z(t^* - \hat{m}_d) - 1 \right) \right. \\ & - Z'(t_{mix}^* - \hat{m}_d) \frac{\partial t_{mix}^*}{\partial m'_d} \left( (1 - Z(t_{mix}^* - \hat{m}_m))K + rZ(t_{mix}^* - m'_d) + (1-r)Z(t_{mix}^* - \hat{m}_d) - 1 \right) \\ & \left. + \left( 1 - Z(t_{mix}^* - \hat{m}_d) \right) \left( -Z'(t_{mix}^* - \hat{m}_m)K \frac{\partial t_{mix}^*}{\partial m'_d} + rZ'(t_{mix}^* - m'_d) \left( \frac{\partial t_{mix}^*}{\partial m'_d} - 1 \right) + (1-r)Z'(t_{mix}^* - \hat{m}_d) \frac{\partial t_{mix}^*}{\partial m'_d} \right) \right]_{m'_d = \hat{m}_d}. \end{aligned}$$

By first making the substitutions  $t_{mix}^*(\hat{m}_d, \hat{m}_d, \hat{m}_m) = t^*(\hat{m}_d, \hat{m}_m)$  and

$$Z'(t^* - \hat{m}_m) = Z'(t^* - \hat{m}_d)/K, \text{ and then } \left. \frac{\partial t_{mix}^*}{\partial m'_d} \right|_{m'_d = \hat{m}_d} = r \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) \text{ and } r = f\rho, \text{ we can write}$$

and simplify this as:

$$\begin{aligned} & -E_{\rho, f} \left[ \rho Z'(t^* - \hat{m}_d) \left( (1 - Z(t^* - \hat{m}_m))K - (1 - Z(t^* - \hat{m}_d)) \right) \right. \\ & \left. - Z'(t^* - \hat{m}_d) \left( \left. \frac{\partial t_{mix}^*}{\partial m'_d} \right|_{m'_d = \hat{m}_d} (1 - Z(t^* - \hat{m}_m))K + \left. \left( r - \frac{\partial t_{mix}^*}{\partial m'_d} \right) \right|_{m'_d = \hat{m}_d} (1 - Z(t^* - \hat{m}_d)) \right) \right] \\ & = E_{\rho, f} \left[ Z'(t^* - \hat{m}_d) \left( \left( \rho f \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) - \rho \right) (1 - Z(t^* - \hat{m}_m))K + \left( \rho f - \rho f \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) + \rho \right) (1 - Z(t^* - \hat{m}_d)) \right) \right] \\ & = Z'(t^* - \hat{m}_d) E_{\rho, f} \left[ \rho \left( \left( f \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) - 1 \right) (1 - Z(t^* - \hat{m}_m))K + \left( f \left( \frac{1}{2} + \frac{\ln K}{d^2} \right) + 1 \right) (1 - Z(t^* - \hat{m}_d)) \right) \right]. \end{aligned}$$

The fitness gradient has the same sign as the expectation, which can be written:

$$\begin{aligned} & E_{\rho, f} [\rho f] \left( \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) (1 - Z(t^* - \hat{m}_m))K + \left( \frac{1}{2} + \frac{\ln K}{d^2} \right) (1 - Z(t^* - \hat{m}_d)) \right) + \\ & E_\rho [\rho] (1 - Z(t^* - \hat{m}_d) - (1 - Z(t^* - \hat{m}_m))K) \\ & = \bar{r} \left( \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) (1 - Z(t^* - \hat{m}_m))K + \left( \frac{1}{2} + \frac{\ln K}{d^2} \right) (1 - Z(t^* - \hat{m}_d)) \right) \\ & + \bar{\rho} ((1 - Z(t^* - \hat{m}_d)) - (1 - Z(t^* - \hat{m}_m))K). \end{aligned}$$

Assuming  $\rho$  and  $f$  are uncorrelated, this equals

$$\bar{\rho} \left( \left( \bar{f} \left( \frac{1}{2} - \frac{\ln K}{d^2} \right) - 1 \right) \left( 1 - Z \left( -\frac{d}{2} - \frac{\ln K}{d} \right) \right) K - \left( \bar{f} \left( -\frac{1}{2} - \frac{\ln K}{d^2} \right) - 1 \right) \left( 1 - Z \left( \frac{d}{2} - \frac{\ln K}{d} \right) \right) \right).$$