

Supporting Information

Eaton and Hallett 10.1073/pnas.1323007111

SI Appendix

The supplementary text provides a detailed description of the mathematical model. *SI1 Mathematical Model* gives a non-technical description of the model. *SI2 Model Equations* gives a mathematical description of the model. *SI3 Model Calibration* describes the calibration of the model to the South African epidemiologic data.

SI1 Mathematical Model

The model uses ordinary differential equations to simulate heterosexual HIV transmission in a two-sex population. The sexually active population is divided into three sexual risk groups that mix semiassortatively. Individuals may move among risk groups, and sexual behavior may change over the course of the epidemic. The model includes progression through five stages of HIV infection according to CD4 cell count, and allows for initiation of ART at any level of CD4 cell count.

1.1 Population Structure. The model simulates a two-sex adult population aged 15 and older (see diagram in Fig. S6). The population is divided into the age groups 15–49, who are presumed to be sexually active, and 50 and older, who are assumed to not form new sexual contacts. Individuals move from the younger to older age group at a rate $\nu = 1/35$ per year and die from the age 50+ population at an annual rate $\mu = 1/11.45$, calibrated to match the relative sizes of the 15–49-y-old population and the age 50+ population in 1990 (1). Individuals enter the age 15–49 population at a rate $\alpha = 0.0226$ per year, calibrated such that the population growth over the period 1990–2010 approximately matches the population growth estimates over that period published by Statistics South Africa (1). All new sexual contacts, and hence HIV transmission, occur in the 15–49 age group, but the older age group is included in the model to assess the total ART need in the intervention scenarios. The sexually active population is divided into three sexual risk groups (termed “low,” “medium,” and “high”). As a crude means of simulating natural variability in individuals’ sexual risk behavior, individuals move from high risk to medium risk, high to low risk, and medium risk to low risk at a rate ψ . This rate is varied in the model calibration (*SI3 Model Calibration*), and is the same for both sexes.

The proportion of the population in each risk group before the introduction of HIV into the population is allowed to be different for each sex and is estimated in the model calibration. The proportion of new entrants into the population who enter each risk group is calculated such that, in the absence of HIV, this proportion of the population in each risk group would remain constant. The proportion entering each risk group remains fixed over the duration of the simulation (meaning that the relative size of risk groups may change as a result of the differential burden of HIV in each risk group).

1.2 Sexual Mixing. As described in the previous section, the population is divided into three sexual risk groups (Fig. S6). The sexual contact rate in each of the risk groups is determined by the overall population average sexual contact rate, the relative rate of new sexual contacts among the three risk groups, and the size of the risk group. The underlying population average rate of new unprotected sexual contacts, $\bar{c}(t)$, is allowed to vary over the course of the epidemic to model potential behavior change in response to the epidemic, such as increased condom use (2) or reductions in the number of new sexual partners (3). The functional form for the reduction is a logistic function parameterized by the initial contact

rate, \bar{c}_0 ; the proportion reduction in the contact rate that will occur, $\Delta\bar{c}$; the start year of the behavior change, $t_{\bar{c}}$; and the number duration (in years) over which the behavior change occurs, $d_{\bar{c}}$. Together, the overall average contact rate at time t is given by

$$\bar{c}(t) = \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot \frac{1}{1 + \exp\left(\frac{t - (t_{\bar{c}} + d_{\bar{c}}/2)}{d_{\bar{c}}/10}\right)} \quad [\text{S1}]$$

The form of this logistic function ensures that the modeled behavior change is symmetric around the midpoint of the behavior change period $t_{\bar{c}} + d_{\bar{c}}/2$ and that the change in behavior occurs in the interval $t_{\bar{c}}$ and $t_{\bar{c}} + d_{\bar{c}}$. To see this, observe that $\bar{c}(t_{\bar{c}}) \approx \bar{c}_0$ and $\bar{c}(t_{\bar{c}} + d_{\bar{c}}) \approx \bar{c}_0 \cdot (1 - \Delta\bar{c})$:

$$\bar{c}(t_{\bar{c}}) = \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot \frac{1}{1 + \exp\left(\frac{t_{\bar{c}} - (t_{\bar{c}} + d_{\bar{c}}/2)}{d_{\bar{c}}/10}\right)} \quad [\text{S2}]$$

$$= \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot \frac{1}{1 + \exp\left(\frac{-d_{\bar{c}}/2}{d_{\bar{c}}/10}\right)} \quad [\text{S3}]$$

$$= \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot \frac{1}{1 + \exp(-5)} \quad [\text{S4}]$$

$$= \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot 0.993 \quad [\text{S5}]$$

$$\approx \bar{c}_0 \quad [\text{S6}]$$

and

$$\bar{c}(t_{\bar{c}} + d_{\bar{c}}) = \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot \frac{1}{1 + \exp\left(\frac{t_{\bar{c}} + d_{\bar{c}} - (t_{\bar{c}} + d_{\bar{c}}/2)}{d_{\bar{c}}/10}\right)} \quad [\text{S7}]$$

$$= \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot \frac{1}{1 + \exp\left(\frac{d_{\bar{c}}/2}{d_{\bar{c}}/10}\right)} \quad [\text{S8}]$$

$$= \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot \frac{1}{1 + \exp(5)} \quad [\text{S9}]$$

$$= \bar{c}_0 \cdot (1 - \Delta\bar{c}) + \bar{c}_0 \cdot \Delta\bar{c} \cdot 0.007 \quad [\text{S10}]$$

$$\approx \bar{c}_0 \cdot (1 - \Delta\bar{c}). \quad [\text{S11}]$$

This is illustrated in the diagram in Fig. S7.

Each of these parameters (the initial contact rate, the percentage reduction in the contact rate, the timing of the start of the reduction, and the duration of the change) is estimated in the model calibration. The relative contact rates between high- and low-risk females and medium- and low-risk females also are estimated. The relative contact rates for males are calculated based on the relative contact rates for females and the sizes of each of

the risk groups for males and females so that the total number of contacts offered by males and females in the same risk group is the same.

The number of sexual contacts formed between members of each risk group is determined by the sexual mixing parameter ε , as proposed by Garnett and Anderson (4). A proportion ε (between 0 and 1) of sexual contacts is reserved exclusively to be formed with other members of the same risk group, whereas the remaining $(1 - \varepsilon)$ proportion of partnerships is formed at random. The value of ε is estimated in the model calibration. If $\varepsilon = 0$, sexual mixing is completely random, whereas if $\varepsilon = 1$, mixing is fully assortative. As HIV mortality differentially affects males and females of each risk group, the total number of contacts offered by males may not be equal to the number of partnerships offered by females. In this case, the number of contacts desired by males and females is geometrically weighted by the parameter θ_G ranging between zero and one; $\theta_G = 0$ indicates that the females' preferences determine the number of contacts, whereas $\theta_G = 1$ indicates that the males' desired number of contacts dominates. For this exercise, the value is fixed at $\theta_G = 0.5$.

1.3 Natural History of HIV Infection. HIV infection is divided into stages according to the CD4 cell count progression associated with the duration of infection (Fig. S8). The stages are:

- i) Primary infection
- ii) CD4 count greater than 350 cells/ μ L
- iii) CD4 count between 200 and 350 cells/ μ L
- iv) CD4 count between 100 and 200 cells/ μ L and
- v) CD4 count below 100 cells/ μ L

HIV-infected individuals progress from one stage of HIV infection to the next at a rate that is the reciprocal of the average duration in the stage. The average duration of primary infection is 2.9 mo, as estimated by Hollingsworth et al. (5). The rates of progression to subsequent CD4 cell stages and the overall duration from HIV infection to death are based on estimates of the time to CD4 cell count thresholds in the sub-Saharan African cohort from the eART-linc (eligibility for ART in lower income countries) collaboration (6), which estimated mean durations of 4.8 y and 9.4 y to reaching CD4 cell counts below 350 and 200 cells/ μ L, respectively. The estimate of an average of 4.17 y between CD4 count ≤ 200 cells/ μ L and CD4 ≤ 100 cells/ μ L is based on an extrapolation of the rate of decline in square root-transformed CD4 cell count between CD4 ≤ 350 cells/ μ L and CD4 ≤ 200 . The overall average duration from infection to HIV death is 14.6 y, from the eART-linc collaboration (6), and the resulting median duration from infection to HIV death is 13.2 y.

The HIV transmission rate of an infected individual varies according to these stages of infection. The overall baseline weighted average transmission rate over the period from the end of primary HIV infection to 1.6 y before death is set to be 0.106 per year, as estimated by Hollingsworth et al. (5) using data from discordant couples in Rakai, Uganda (7). The relative rates of transmission during the CD4 stages CD4 > 350 , CD4 200–350, and CD4 ≤ 200 are based on the relative rates of transmission observed for these stages by Donnell et al. (8), although the rate of transmission during the CD4 ≤ 100 stages has been reduced in accordance with the estimate from Hollingsworth et al. (5) that no transmission occurs during the final 9 mo of infection, presumably because individuals are sick and not very sexually active during this period. The transmission rate during primary HIV infection is set at 2.76 per year, as estimated by Hollingsworth et al. (5).

1.4 ART Model. Individuals may initiate ART from any of the above stages of HIV infection. ART is divided into a multistage process (Fig. S9). Upon treatment initiation, all individuals first enter a “virally suppressing” stage during which they are on ART but their viral load is not yet fully suppressed. This stage lasts for a

mean of 3 mo, such that 86% of patients achieve virological suppression after 6 mo, consistent with levels of viral suppression after 6 mo in ART-naïve patients in the United Kingdom (9). Transmission is assumed to be reduced by half during this stage compared with the CD4 stage from which treatment was initiated.

After this stage, most patients enter a long period of “effective ART,” whereas a proportion of patients for whom treatment is not successful go directly to the final stage of the ART model of being “very sick,” which lasts for an average of 6.2 mo before death. The probability of immediately failing treatment depends on the CD4 cell count stage from which treatment was initiated to allow for high early mortality when starting treatment at lower CD4 cell counts, but then relatively similar long-term survival if treatment effectively suppresses viral load and symptoms are controlled (10–12). The proportion of patients who fail ART is calibrated such that mortality in the first year after initiating ART matches the crude first-year mortality rate observed for each CD4 count stratum in a collaborative analysis of ART cohorts from sub-Saharan Africa (11).

For most patients in whom ART is effective, the viral load is suppressed and transmission is reduced by 92% compared with the HIV transmission rate in the CD4 200–350 cells/ μ L stage (8). The period of effective ART is divided into two stages—first, a period of early effective ART lasting an average of 1.75 y, and then a long period of sustained viral suppression. The reduction in transmission is assumed to be the same in both these stages, but this is implemented as separate ART stages so that the dropout rate from treatment can be varied according to the duration on ART. For example, we may assume that there is high dropout in the years following ART initiation, but patients who remain on treatment for 2 y likely have accommodated treatment and have high retention thereafter. In addition to the previously described higher probability for immediate treatment failure and death for those starting ART at low CD4 cell counts, the failure rate for long-term effective ART is assumed to be slightly lower for those who start treatment at high CD4 cell counts (Fig. S9), in line with observations that mortality is modestly higher even several years after treatment initiation for those who start at low CD4 cell counts (12) and to ensure there is no “survival benefit” in the model from delaying treatment initiation.

After patients fail long-term effective ART, they enter a stage of “treatment failing” in which they are viremic and are assumed to have the same infectiousness as individuals in the CD4 cell count category 200–350 cells/ μ L. The average duration of this stage is 2.3 y. Finally, individuals enter a stage of being very sick just before death, which lasts an average of 6.2 mo. During this period of being very sick, transmission is reduced and assumed to be at the same level as during the CD4 ≤ 100 cells/ μ L stage.

1.5 Dropping out from ART. Individuals may drop out from any of the first three stages of ART: virally suppressing, early effective ART, and effective ART. The rate of dropout from treatment is permitted to vary according to duration on ART and the CD4 cell count category from which treatment was initiated, in line with data suggesting that those starting treatment at higher CD4 cell counts may have poorer retention in treatment programs (13, 14). Rates of dropout from ART are representative of those observed in South African ART cohorts (12, 15–18). Assumed rates of dropping out from ART are presented in Table S2. Dropout is assumed to decline after the first 2 y on treatment, and individuals who reinstate treatment after having dropped out once (see below) are assumed to have lower rates of dropout.

After dropping out from treatment, the untreated CD4 stage category that individuals enter depends on their pretreatment CD4 category and the duration on treatment to simulate CD4 cell count reconstitution associated with ART. Individuals who drop out from the virally suppressing stage all return to the same CD4 stage from which they initiated treatment. For those who drop out

during the early effective ART stage, half move one CD4 stage higher, whereas half increase two CD4 stages. Those who drop out from the “effective ART” stage all increase two CD4 stages. This is summarized in Table S3. However, individuals who have dropped out of treatment progress through subsequent CD4 stages twice as fast as treatment-naïve individuals (rates described in Fig. S9).

In the model, after individuals drop out of treatment, they are eligible to restart treatment again once. The rate at which these individuals reinitiate treatment depends on the CD4 cell count category, such that they are increasingly likely to reinitiate treatment at lower CD4 count categories, when they likely are experiencing clinical symptoms. The rate of reinitiating treatment for those with CD4 cell counts between 200 and 350 cells/ μL is 0.048 per year, the rate for those with CD4 cell counts between 100 and 200 cells/ μL is 0.160, and the rate for those with CD4 counts below 100 cells/ μL is 2.92 per year. Based on progression through these stages at twice the rate of treatment-naïve individuals, the probability that an individual who dropped out will restart ART in each of these CD4 cell count categories before progressing to the next CD4 stage (or dying, in the case of those with CD4 counts ≤ 100 cells/ μL) is 10%, 25%, and 60%, respectively. The overall probability that an individual who has dropped out of treatment once will reinitiate treatment before dying from HIV is 73%.

Upon restarting treatment, individuals progress through the same stages of ART as when first initiating ART. The baseline assumption for the rate of dropping out from treatment after reinitiating is 0.06 per year (Table S2, last column). Because individuals may restart treatment, but only once, the model separately tracks treated and untreated people according to the number of times they have initiated ART. To summarize this, Table S1 lists all the stages of antiretroviral treatment through which infected individuals may progress, and the subscript identifying each stage in the technical model description that follows.

1.6 HIV Transmission. The probability of transmission during contact between a susceptible and an infected individual depends on the annual transmission rate $\beta_{m,u}$ in HIV stage m and treatment stage u , the “intensity” of a contact κ_{r_M, r_F} between a male in risk group r_M and a female in risk group r_F . The intensity parameter accounts for factors that affect the probability of transmission in different types of partnerships, such as the partnership duration, coital frequency, and condom use. The per-contact transmission probability based on these parameters is $1 - e^{-\beta_m \kappa_{r_M, r_F}}$, where m is the HIV stage of the infected partner.

1.7 Epidemic Seeding. We initialize the HIV epidemic with an adult HIV prevalence of 0.025% at time t_0 (estimated). This initial prevalence is distributed across the sexes, risk groups, and infected stages proportional to the eigenvector associated with the largest eigenvalue of the linearized Jacobian matrix describing transmission in a fully susceptible population (the matrix $\mathbf{T} + \Sigma$ defined in *SI4 R₀ Calculation*). Thus, the initial distribution of infections is consistent with that which would be expected during the early exponential growth period for a given set of parameters.

SI2 Model Equations

We divide the population into four categories according to their infection and treatment status:

$S^{g,r}$: HIV-uninfected and sexually active (susceptible) individuals of sex g in risk group r .

$I_{m,u}^{g,r}$: HIV-infected and sexually active individuals of sex g in risk group r and HIV infection stage m . The subscript $u=0$ indicates untreated individuals who do not have access to

treatment, $u=6$ indicates individuals who have dropped out of treatment and are eligible to restart, and $u=12$ indicates those who have dropped out from treatment a second time and are not eligible to restart.

$T_{m,u}^{g,r}$: HIV-infected and sexually active individuals on ART (treated) who began treatment in HIV infection stage m and are in treatment stage u (Table S1).

$R_{m,u}^g$: Individuals removed from the sexually active population of sex g , in HIV infection stage m , and treatment stage u . (Uninfected individuals are indicated by $m=0$ and untreated individuals by $u \in \{0, 6, 12\}$.)

In the above and throughout the following mathematical description, the superscript $g \in \{M, F\}$ corresponds to sex; superscript $r \in \{H, M, L\}$ corresponds to the sexual risk group for the sexually active population; subscript $m \in \{0, \dots, 5\}$ corresponds to HIV infection stage, with 0 representing uninfected and 1–5 corresponding to the stages of infection from primary infection to CD4 ≤ 100 ; and subscript $u \in \{0, \dots, 12\}$ corresponds to ART status, with 0 indicating untreated individuals without access to treatment, and the remaining stages indicating different stages of ART as indicated in Table S1.

The following differential equations define the dynamics of the groups:

$$\frac{dS^{g,r}}{dt} = \frac{\alpha + \nu}{2} \tilde{\pi}^{g,r} (S^{\cdot,\cdot} + I^{\cdot,\cdot} + T^{\cdot,\cdot}) + \sum_{r'} \Psi_{r',r} \cdot S^{g,r'} - \left(f^{g,r}(t) + \nu + \sum_{r'} \Psi_{r,r'}^g \right) S^{g,r}$$

$$\frac{dI_{1,0}^{g,r}}{dt} = f^{g,r}(t) S^{g,r} + \sum_{r'} \Psi_{r',r} \cdot I_{1,0}^{g,r'} - \left(\sigma_1 + \nu + \sum_{r'} \Psi_{r,r'}^g \right) I_{1,0}^{g,r}$$

$$\frac{dI_{m,0}^{g,r}}{dt} = \sigma_{m-1} I_{m-1,0}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot I_{m,0}^{g,r'} - \left(\sigma_m + \nu + \sum_{r'} \Psi_{r,r'}^g \right) I_{m,0}^{g,r}$$

for $m \geq 2$

$$\frac{dT_{m,1}^{g,r}}{dt} = \lambda_m^g I_{m,0}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot T_{m,1}^{g,r'} - \left(\phi_{m,1}^g + \eta_{m,1} + \nu + \sum_{r'} \Psi_{r,r'}^g \right) T_{m,1}^{g,r}$$

$$\frac{dT_{m,2}^{g,r}}{dt} = (1 - \xi_m) \phi_{m,1}^g T_{m,1}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot T_{m,2}^{g,r'} - \left(\phi_{m,2}^g + \eta_{m,2} + \nu + \sum_{r'} \Psi_{r,r'}^g \right) T_{m,2}^{g,r}$$

$$\frac{dT_{m,u}^{g,r}}{dt} = \phi_{m,u-1}^g T_{m,u-1}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot T_{m,u}^{g,r'} - \left(\phi_{m,u}^g + \eta_{m,u} + \nu + \sum_{r'} \Psi_{r,r'}^g \right) T_{m,u}^{g,r} \quad \text{for } u \in \{3, 4\}$$

$$\frac{dT_{m,5}^{g,r}}{dt} = \xi \phi_{m,1}^g T_{m,1}^{g,r} + \phi_{m,4}^g T_{m,4}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot T_{m,5}^{g,r'} - \left(\phi_{m,5}^g + \eta_{m,5} + \nu + \sum_{r'} \Psi_{r,r'}^g \right) T_{m,5}^{g,r}$$

$$\frac{dI_{m,6}^{g,r}}{dt} = \sum_{m'} \sum_{u=1}^3 \rho_{m',u}^m \eta_{m',u} T_{m',u}^{g,r} + \tilde{\sigma}_{m-1} I_{m-1,6}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot I_{m,6}^{g,r'} - \left(\tilde{\sigma}_m + \tilde{\lambda}_m^g + \nu + \sum_{r'} \Psi_{r',r}^g \right) I_{m,6}^{g,r} \quad \text{for } m \geq 1$$

$$\frac{dT_{m,7}^{g,r}}{dt} = \tilde{\lambda}_m^g I_{m,6}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot T_{m,7}^{g,r'} - \left(\Phi_{m,7}^g + \eta_{m,7} + \nu + \sum_{r'} \Psi_{r',r}^g \right) T_{m,7}^{g,r}$$

$$\frac{dT_{m,8}^{g,r}}{dt} = (1 - \xi_m) \Phi_{m,7}^g T_{m,7}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot T_{m,8}^{g,r'} - \left(\Phi_{m,8}^g + \eta_{m,8} + \nu + \sum_{r'} \Psi_{r',r}^g \right) T_{m,8}^{g,r}$$

$$\frac{dT_{m,u}^{g,r}}{dt} = \Phi_{m,u-1}^g T_{m,u-1}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot T_{m,u}^{g,r'} - \left(\Phi_{m,u}^g + \eta_{m,u} + \nu + \sum_{r'} \Psi_{r',r}^g \right) T_{m,u}^{g,r} \quad \text{for } u \in \{9, 10\}$$

$$\frac{dT_{m,11}^{g,r}}{dt} = \xi \Phi_{m,7}^g T_{m,7}^{g,r} + \Phi_{m,10}^g T_{m,10}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot T_{m,11}^{g,r'} - \left(\Phi_{m,11}^g + \eta_{m,11} + \nu + \sum_{r'} \Psi_{r',r}^g \right) T_{m,11}^{g,r}$$

$$\frac{dT_{m,12}^{g,r}}{dt} = \sum_{m'} \sum_{u=6}^9 \rho_{m',u}^m \tilde{\eta}_{m',u} T_{m',u}^{g,r} + \tilde{\sigma}_{m-1} I_{m-1,12}^{g,r} + \sum_{r'} \Psi_{r',r} \cdot I_{m,12}^{g,r'} - \left(\tilde{\sigma}_m + \tilde{\lambda}_m^g + \nu + \sum_{r'} \Psi_{r',r}^g \right) I_{m,12}^{g,r} \quad \text{for } m \geq 1$$

$$\frac{dR_{0,0}^g}{dt} = \nu \sum_r S^{g,r} - \mu R_{0,0}^g$$

$$\frac{dR_{m,0}^g}{dt} = \sigma_{m-1} R_{m-1,0}^g + \nu \sum_r I_m^{g,r} - (\sigma_m + \mu) R_{m,0}^g \quad \text{for } m \geq 1$$

$$\frac{dR_{m,1}^{g,r}}{dt} = \lambda_m^g R_{m,0}^{g,r} + \nu \sum_r T_{m,1}^{g,r} - \left(\Phi_{m,1}^g + \eta_{m,1} + \mu \right) R_{m,1}^{g,r}$$

$$\frac{dR_{m,2}^{g,r}}{dt} = (1 - \xi_m) \Phi_{m,1}^g R_{m,1}^{g,r} + \nu \sum_r T_{m,2}^{g,r} - \left(\Phi_{m,2}^g + \eta_{m,2} + \mu \right) R_{m,2}^{g,r}$$

$$\frac{dR_{m,u}^{g,r}}{dt} = \Phi_{m,u-1}^g R_{m,u-1}^{g,r} + \nu \sum_r R_{m,u}^{g,r} - \left(\Phi_{m,u}^g + \eta_{m,u} + \mu \right) R_{m,u}^{g,r} \quad \text{for } u \in \{3, 4\}$$

$$\frac{dR_{m,5}^{g,r}}{dt} = \xi \Phi_{m,1}^g R_{m,1}^{g,r} + \Phi_{m,4}^g R_{m,4}^{g,r} + \nu \sum_r T_{m,5}^{g,r} - \left(\Phi_{m,5}^g + \eta_{m,5} + \mu \right) R_{m,5}^{g,r}$$

$$\frac{dR_{m,6}^{g,r}}{dt} = \sum_{m'} \sum_{u=1}^3 \rho_{m',u}^m \eta_{m',u} R_{m',u}^{g,r} + \tilde{\sigma}_{m-1} R_{m-1,6}^{g,r} + \nu \sum_r I_{m,6}^{g,r} - \left(\tilde{\sigma}_m + \tilde{\lambda}_m^g + \mu \right) I_{m,6}^{g,r} \quad \text{for } m \geq 1$$

$$\frac{dR_{m,7}^{g,r}}{dt} = \tilde{\lambda}_m^g R_{m,6}^{g,r} + \nu \sum_r T_{m,7}^{g,r} - \left(\Phi_{m,7}^g + \eta_{m,7} + \mu \right) R_{m,7}^{g,r}$$

$$\frac{dR_{m,8}^{g,r}}{dt} = (1 - \xi_m) \Phi_{m,7}^g R_{m,7}^{g,r} + \nu \sum_r T_{m,8}^{g,r} - \left(\Phi_{m,8}^g + \eta_{m,8} + \mu \right) R_{m,8}^{g,r}$$

$$\frac{dR_{m,u}^{g,r}}{dt} = \Phi_{m,u-1}^g R_{m,u-1}^{g,r} + \nu \sum_r T_{m,u}^{g,r} - \left(\Phi_{m,u}^g + \eta_{m,u} + \mu \right) R_{m,u}^{g,r} \quad \text{for } u \in \{9, 10\}$$

$$\frac{dR_{m,11}^{g,r}}{dt} = \xi \Phi_{m,7}^g R_{m,7}^{g,r} + \Phi_{m,10}^g R_{m,10}^{g,r} + \nu \sum_r T_{m,11}^{g,r} - \left(\Phi_{m,11}^g + \eta_{m,11} + \mu \right) R_{m,11}^{g,r}$$

$$\frac{dR_{m,12}^{g,r}}{dt} = \sum_{m'} \sum_{u=7}^9 \rho_{m',u}^m \tilde{\eta}_{m',u} R_{m',u}^{g,r} + \tilde{\sigma}_{m-1} R_{m-1,12}^{g,r} + \nu \sum_r I_{m,12}^{g,r} - \left(\tilde{\sigma}_m + \tilde{\lambda}_m^g + \mu \right) R_{m,12}^{g,r} \quad \text{for } m \geq 1.$$

In the above equations, the parameter $\tilde{\pi}^{g,r}$ is the proportion of new susceptible individuals of sex g that should enter risk group r to maintain a constant proportion $\pi^{g,r}$ in risk group r in the absence of HIV infection. Solving the above equations with $f^{g,r}(t) = 0$ gives that

$$\tilde{\pi}^{g,r} = \frac{\sum_{r'} \Psi_{r',r}^g \pi^{g,r'} - \sum_{r'} \Psi_{r',r}^g \pi^{g,r}}{\alpha + \nu}. \quad \text{[S12]}$$

The function $f^{g,r}(t)$ is the force of infection for the group $S^{g,r}$ that depends on the contact rate $c^{g,r}(t)$ at time t in that group, the probability that a contact is formed with an infectious partner, and the probability of transmission in that contact.

The contact rate for a risk group depends on the average underlying contact rate $\bar{c}(t)$, which changes over time according to Eq. S1, the size of the risk groups at the beginning of the epidemic $\pi^{g,r}$, and the relative contact rates of the high- and medium-risk groups to the low-risk group $\omega^{g,r}$ (where $\omega^{g,L} := 1$). The weighted average of the relative contact rate by the size of the initial risk group yields the annual contact rate for each risk group at time t :

$$c^{g,r}(t) = \frac{\bar{c}(t) \cdot \omega^{g,r}}{\sum_r \pi^{g,r} \omega^{g,r}}.$$

The total number of contacts desired to be formed by members of risk group r of sex g thus is

$$J^{g,r}(t) = c^{g,r}(t) \left(S^{g,r} + \sum_m I_m^{g,r} + \sum_m \sum_u T_{m,u}^{g,r} \right). \quad \text{[S13]}$$

The number of these contacts $J^{g,r}$ desired to be formed with each risk group r' of the opposite sex g' depends on the value of

the assortativity parameter ε . A proportion ε of the partnerships are desired to be formed only with the members of the same risk group $r=r'$, whereas the remaining $1-\varepsilon$ proportion of the partnerships are formed with each risk group of the opposite sex proportionally to the number of partnerships $J^{g,r'}$ offered by those risk groups. Formally, the proportion of contacts desired to be formed by sex g and risk group r that are formed with the risk group r' of the opposite sex is defined as

$$Q_{r,r'}^g = \varepsilon \delta_{r,r'} + (1-\varepsilon) \frac{J^{g,r'}}{\sum_r J^{g,r'}}, \quad [\text{S14}]$$

where $\delta_{r,r'}$ is the Kronecker delta defined as $\delta_{r,r'} = 1$ if $r=r'$ and 0 otherwise.

In the case that males in risk group r_M and females in group r_F do not agree on the number of partnerships to be formed between the risk groups, i.e., $Q_{r_M,r_F}^M J^{M,r_M} \neq Q_{r_F,r_M}^F J^{F,r_F}$, the discrepancy is balanced according to the parameter θ_G as

$$\begin{aligned} \tilde{Q}_{r_M,r_F}^M &= Q_{r_M,r_F}^M \left(\frac{Q_{r_M,r_F}^M J^{M,r_M}}{Q_{r_F,r_M}^F J^{F,r_F}} \right)^{\theta_G - 1} \\ \tilde{Q}_{r_F,r_M}^F &= Q_{r_F,r_M}^F \left(\frac{Q_{r_M,r_F}^M J^{M,r_M}}{Q_{r_F,r_M}^F J^{F,r_F}} \right)^{\theta_G}. \end{aligned} \quad [\text{S15}]$$

The probability that transmission occurs in a contact between a susceptible and an infected individual depends on the stage of infection and treatment status of the infection according to the transmission rate parameter $\beta_{m,u}$ and on the value of the partnership intensity multiplier κ_{r_M,r_F} for a partnership between the male in risk group r_M and the female in risk group r_F . The force of infection then is calculated by summing over the probability that each contact is with an infectious individual and the probability that transmission occurs according to the infection stage and treatment status of the partner:

$$f^{g,r}(t) = c^{g,r}(t) \sum_{r'} \left[\tilde{Q}_{r,r'}^g \frac{\sum_m \sum_u J_m^{g,r'} P_{m,u}^{g,r',r} + \sum_m \sum_u T_m^{g,r'} P_{m,u}^{g,r',r}}{S^{g,r'} + \sum_m \sum_u I_m^{g,r} + \sum_m \sum_u T_m^{g,r}} \right], \quad [\text{S16}]$$

where $P_{m,u}^{g,r'}$ is the probability of transmission per contact by an infected individual of sex g in risk group r , HIV stage m , and treatment status u to a susceptible individual of the opposite sex in risk group r' , defined as

$$P_{m,u}^{g,r'} = 1 - \exp\{-\beta_{m,u} \cdot \kappa_{r_M,r_F}\}. \quad [\text{S17}]$$

S13 Model Calibration

The model is calibrated to nationally representative HIV prevalence data and ART scale-up data from South Africa. The general strategy for model calibration is that parameters related to the natural history of HIV infection, the effect of ART on individual infection, and patterns of access and retention in the existing ART program are fixed and informed from the literature as described in *S11 Mathematical Model*. Parameters relating to sexual behavior and mixing, the start time of the epidemic, the timing and magnitude of sexual behavior change, and the timing and rate of existing ART scale-up are estimated using a Bayesian approach. This yields a joint distribution of parameter combinations representing different sexual mixing patterns consistent with the observed epidemic.

3.1 Data. The model is calibrated using HIV prevalence data from two sources. The first is HIV prevalence among 15–49-y-old

pregnant women from the annual national antenatal prevalence surveys from 1990 to 2008 (19). The second is national HIV prevalence among 15–49-y-old males and females from the three nationally representative household surveys conducted by the Human Sciences Research Council in 2002, 2005, and 2008 (20–22). The discrepancy between the level of HIV prevalence in the antenatal surveillance and the household survey-based prevalence is reconciled by incorporating a bias parameter in the antenatal prevalence compared with prevalence among the general 15–49-y-old population from the household surveys. This is described in detail in section 3.3.

The model also is calibrated to the percentage of the adult population on ART. This is calculated by dividing the total number of adults reported to be on ART according to the South Africa Department of Health in June of each year from 2005 to 2010 (23) by the annual midyear population size estimate of those over 15 y old from Statistics South Africa (1). The resulting estimates for the proportion of the adult population on ART for 2005–2010 to which the model is calibrated are shown in Fig. S10.

3.2 Estimated Model Parameters. A vector θ of 17 model parameters are estimated, and are used either directly or to derive a number of the model inputs in the equations described in *S12 Model Equations* and Table S4. The mathematical model parameters that are estimated from fitting to HIV prevalence data from South Africa are given in Table S5, along with the prior distribution from which they are estimated. This collection of model parameters estimated from the data will be referenced together as the parameter vector θ .

3.3 Statistical Methods. We use a Bayesian analysis to estimate probability distributions for the unknown parameters given the model and available HIV prevalence data, \mathcal{W} . Let θ denote the set of parameters to be estimated from section 3.2 and M denote the mathematical model and the associated fixed parameter values. For a given set of parameter values, the model produces a corresponding set of predicted HIV prevalence values $\zeta = M(\theta)$. Then we specify a likelihood function $p(\mathcal{W}|M(\theta))$ for the probability of the data given the model and parameter values. If we let $p(\theta)$ denote a prior distribution on the unknown parameters, using Bayes theorem the posterior distribution of θ given the model and data are given by

$$p(\theta|\mathcal{W}, M) \propto p(\theta)p(\mathcal{W}|M(\theta)).$$

3.3.1 Likelihood function. We now derive the combined likelihood for the national seroprevalence survey data and antenatal clinic data. First, we define the likelihood for an individual datum.

For the national household survey estimates, let $W_{F,t}$ and $W_{M,t}$ be the HIV prevalence in the age group 15–49 reported by a survey in year t . As above, define $\zeta_{g,t}$ to be the predicted population HIV prevalence for sex g at time t by the model $M(\theta)$. We assume that HIV prevalence estimates from national household surveys are an unbiased estimate of the true population prevalence, but we logit-transform the data to stabilize the error variance so that

$$\log\left(\frac{W_{g,t}}{1-W_{g,t}}\right) = \log\left(\frac{\zeta_{g,t}}{1-\zeta_{g,t}}\right) + \epsilon_{g,t}, \quad [\text{S18}]$$

where $\epsilon_{g,t} \sim \text{Normal}(0, \sigma_{g,t}^2)$ and the errors are conditionally independent given ζ . Thus, the likelihood function for an estimate from a national prevalence survey is

$$p(W_{g,t}|M(\theta)) = \frac{1}{\sqrt{2\pi\sigma_{g,t}^2}} \exp\left\{\frac{-1}{2\sigma_{g,t}^2}(\logit(W_{g,t}) - \logit(\zeta_{g,t}))^2\right\}. \quad [\text{S19}]$$

In calculating the likelihood, the value of $\sigma_{g,t}^2$ is replaced by an estimate $\hat{\sigma}_{g,t}^2$ based on the confidence intervals reported by the Human Sciences Research Council (HSRC) survey, which account for the complex sampling design. The confidence intervals in the HSRC reports are on the inverse-logit scale, so the error variances are estimated by

$$\hat{\sigma}_{g,t}^2 = \left(\frac{\logit(CI_{g,t}^{\max}) - \logit(CI_{g,t}^{\min})}{2 \cdot \Phi^{-1}(.975)}\right)^2, \quad [\text{S20}]$$

where Φ^{-1} is the inverse of the standard normal cumulative distribution function.

For the antenatal clinic data, similar to refs. 2 and 24, we assume that HIV prevalence $W_{C,t}$ among antenatal clinic attendees at time t is linearly related to prevalence in the general female age 15–49 population on the logit scale and that this effect is fixed over time. To model this, we introduce an additional parameter γ such that

$$\log\left(\frac{W_{C,t}}{1 - W_{C,t}}\right) = \log\left(\frac{\zeta_{F,t}}{1 - \zeta_{F,t}}\right) + \gamma + \epsilon_{C,t}, \quad [\text{S21}]$$

where the error term $\epsilon_{C,t} \sim \text{Normal}(0, \sigma_{C,t}^2)$ and recalling that $\zeta_{F,t}$ is the HIV prevalence among females aged 15–49 predicted by the model at time t . Thus,

$$p(W_{C,t}|M(\theta)) = \frac{1}{\sqrt{2\pi\sigma_{C,t}^2}} \exp\left\{\frac{-1}{2\sigma_{C,t}^2}(\logit(W_{C,t}) - (\logit(\zeta_{F,t}) + \gamma))^2\right\}. \quad [\text{S22}]$$

Once again, when evaluating the likelihood, we replace $\sigma_{C,t}^2$ with an estimate $\hat{\sigma}_{C,t}^2$ calculated from the confidence intervals published by the South African Department of Health. In the case of the antenatal clinic data, the reported confidence intervals are symmetric, suggesting that they have been estimated on the untransformed scale rather than logit-transformed prevalence. Deriving an estimate of the sampling error variance on the logit scale involves two steps: first estimating the untransformed error variance $\hat{\tau}_{C,t}^2$, and then using the delta method to approximate the variance $\hat{\sigma}_{C,t}^2$ of the logit-transformed distribution, which depends on the estimated antenatal clinic prevalence at time t , $W_{C,t}$. The equations for this are

$$\hat{\tau}_{C,t}^2 = \left(\frac{\logit(CI_{C,t}^{\max}) - \logit(CI_{C,t}^{\min})}{2 \cdot \Phi^{-1}(.975)}\right)^2 \text{ and} \quad [\text{S23}]$$

$$\hat{\sigma}_{C,t}^2 = \frac{\hat{\tau}_{C,t}^2}{W_{C,t}^2(1 - W_{C,t})^2}. \quad [\text{S24}]$$

To arrive at the likelihood for the full data W , we assume that the data points are conditionally independent given the predicted prevalences $\xi = M(\theta)$ and the antenatal clinic bias parameter γ . Defining T_N to be the set of years for which national survey prev-

alence estimates are available and T_C the set of years for which antenatal clinic estimates are available, the full likelihood is

$$\begin{aligned} p(W|M(\theta), \gamma) &= \prod_{t \in T_N} \prod_{g \in \{M,F\}} p(W_{g,t}|M(\theta)) \cdot \prod_{t \in T_C} p(W_{C,t}|M(\theta), \gamma) \\ &= \prod_{t \in T_N} \prod_{g \in \{M,F\}} \frac{1}{\sqrt{2\pi\sigma_{g,t}^2}} \exp\left\{\frac{-1}{2\sigma_{g,t}^2}(\logit(W_{g,t}) - \logit(\zeta_{g,t}))^2\right\} \\ &\quad \times \prod_{t \in T_C} \frac{1}{\sqrt{2\pi\sigma_{C,t}^2}} \exp\left\{\frac{-1}{2\sigma_{C,t}^2}(\logit(W_{C,t}) - (\logit(\zeta_{F,t}) + \gamma))^2\right\}. \end{aligned} \quad [\text{S25}]$$

3.3.2 Priors. The prior distributions for the estimated model parameters θ are given in Table S5. The antenatal bias parameter γ is assumed to have Uniform(0, 1) prior distribution. This amounts to assuming that the odds ratio of ANC prevalence to adult female prevalence is between 1 and 2.7.

3.3.3 Estimating the posterior distribution. Multiplying the likelihood and prior distribution yields the joint posterior distribution of the parameters θ and γ given the model and data up to a scaling constant:

$$p(\theta, \gamma|W, M) \propto p(W|\theta, \gamma, M)p(\theta, \gamma). \quad [\text{S26}]$$

We principally are interested in the values of the model parameters θ . The posterior distribution for $\theta|W, M$ can be calculated by integrating out the parameter γ . Observe that

$$\begin{aligned} p(\theta|W, M) &\propto \int_{\Omega_\gamma} p(W|\theta, \gamma, M)p(\theta, \gamma) d\gamma \\ &= \int_{\Omega_\gamma} p(\theta)p(\gamma) \prod_{t \in T_N} \prod_{g \in \{M,F\}} p(W_{g,t}|M(\theta)) \cdot \prod_{t \in T_C} p(W_{C,t}|M(\theta), \gamma) d\gamma \\ &= \prod_{t \in T_N} \prod_{g \in \{M,F\}} p(W_{g,t}|M(\theta))p(\theta) \cdot \int_{\Omega_\gamma} p(\gamma) \prod_{t \in T_C} p(W_{C,t}|M(\theta), \gamma) d\gamma; \end{aligned} \quad [\text{S27}]$$

so if we can efficiently evaluate the integral $\int_{\Omega_\gamma} p(\gamma) \prod_{t \in T_C} p(W_{C,t}|M(\theta), \gamma) d\gamma$, then we can efficiently estimate the posterior distribution of θ . Before attacking this, let us define three useful quantities:

$$S^2 = \left(\sum_{t \in T_C} \frac{1}{\sigma_{C,t}^2}\right)^{-1} \quad [\text{S28}]$$

$$\bar{D} = S^2 \cdot \sum_{t \in T_C} \frac{W_{C,t} - \zeta_{F,t}}{\sigma_{C,t}^2} \quad [\text{S29}]$$

$$\bar{D}^2 = S^2 \cdot \sum_{t \in T_C} \frac{(W_{C,t} - \zeta_{F,t})^2}{\sigma_{C,t}^2}. \quad [\text{S30}]$$

The first may be thought of as the pooled variance of the ANC prevalence estimates, the second as the precision-weighted mean difference between ANC data prevalence and the model-predicted female prevalence, and the third as the precision-weight mean-squared distance between the ANC data and the predicted female prevalence.

Now, again consider our integral. We will do a bit of rearranging to show that the integral can be evaluated as a normal cumulative

distribution function. For brevity, denote $\dot{W}_t = \text{logit}(W_{C,t})$ and $\zeta_t = \text{logit}(\zeta_{F,t})$, and all sums and products are over the set T_C .

$$\begin{aligned}
& \int_{\Omega_t} p(\gamma) p(W_C | M(\theta), \gamma) d\gamma \\
&= \int_0^1 \prod_t \frac{1}{\sqrt{2\pi\sigma_{C,t}^2}} \exp\left\{ \frac{-1}{2\sigma_{C,t}^2} (\dot{W}_t - (\zeta_t + \gamma))^2 \right\} d\gamma \\
&= \frac{1}{\sqrt{\prod_t 2\pi\sigma_{C,t}^2}} \int_0^1 \exp\left\{ \sum_t \frac{-1}{2\sigma_{C,t}^2} (\gamma - (\dot{W}_t - \zeta_t))^2 \right\} d\gamma \\
&= K \cdot \int_0^1 \exp\left\{ \frac{-1}{2} \sum_t \left(\frac{\gamma^2}{\sigma_{C,t}^2} - 2\gamma \frac{\dot{W}_t - \zeta_t}{\sigma_{C,t}^2} + \frac{(\dot{W}_t - \zeta_t)^2}{\sigma_{C,t}^2} \right) \right\} d\gamma \\
&= K \cdot \int_0^1 \exp\left\{ \frac{-1}{2} \left(\gamma^2 / S^2 - 2\gamma \bar{D} / S^2 + \bar{D}^2 / S^2 \right) \right\} d\gamma \\
&= K \cdot \int_0^1 \exp\left\{ \frac{-1}{2S^2} (\gamma - \bar{D})^2 + \frac{1}{2S^2} (\bar{D}^2 - \bar{D}^2) \right\} d\gamma \\
&= K \cdot e^{(\bar{D}^2 - \bar{D}^2) / (2S^2)} \int_0^1 \frac{1}{\sqrt{2\pi S^2}} e^{-(\gamma - \bar{D})^2 / (2S^2)} d\gamma \\
&= \sqrt{\frac{2\pi S^2}{\prod_t 2\pi\sigma_{C,t}^2}} e^{(\bar{D}^2 - \bar{D}^2) / (2S^2)} \left[\Phi\left(\frac{1 - \bar{D}}{\sqrt{S^2}}\right) - \Phi\left(\frac{0 - \bar{D}}{\sqrt{S^2}}\right) \right].
\end{aligned} \tag{S31}$$

Using this expression, we can efficiently evaluate the posterior density function $p(\theta | W, M)$ up to a constant. We use the incremental mixture importance sampling algorithm to approximate and sample from the posterior distribution (25).

S14 R_0 Calculation

We calculate $R_0(t)$ over the course of the epidemic as the dominant eigenvalue of the next-generation matrix (NGM) following the formalism for compartmental systems described by Diekmann et al. (26). Consider the Jacobian matrix representing the linearization of the infected subsystem (the equations $I_m^{g,r}$ defining the infected persons, excluding the susceptible $S^{g,r}$, treated $T_{m,u}^{g,r}$, and removed $R_{m,u}^g$ stages). We decompose this into the sum of two matrices $\mathbf{T} + \mathbf{\Sigma}$, where \mathbf{T} describes transmissions giving rise to new infected persons and $\mathbf{\Sigma}$ describes all other transitions between infected states. We calculate R_0 as the dominant eigenvalue of the NGM with large domain, $\mathbf{K}_L := -\mathbf{T}\mathbf{\Sigma}^{-1}$ (26).

Both \mathbf{T} and $\mathbf{\Sigma}$ are 30×30 matrices summarizing the rates of transmission and transition between $g=2$ sexes, $m=5$ stages of HIV infection, and $r=3$ risk groups. The transmission matrix T_{ij} gives the rate of new infections in state i created by an infected person in state j . Because the contact rate $c^{g,r}(t)$ varies over time (according to the parameter $\Delta_{\bar{c}}$; Eq. S1), $\mathbf{T}(t)$ is a function of time t . Because all transmission occurs heterosexually, the transmission matrix \mathbf{T} consists of 15×15 submatrices.

$$\mathbf{T}(t) = \begin{bmatrix} 0 & T^{F \rightarrow M} \\ \mathbf{T}^{M \rightarrow F} & 0 \end{bmatrix}. \tag{S32}$$

All newly infected persons start in the first stage of infection (early transmission), and so the sex-specific transmission matrices con-

sist of a row of 3×3 submatrices $T_m^{g,g'}$ for the rate of transmission from infected persons in stage m with zeros below:

$$\mathbf{T}^{g \rightarrow g'}(t) = \begin{bmatrix} T_1^g & T_2^g & T_3^g & T_4^g & T_5^g \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \tag{S33}$$

$$T_m^g(t) = \begin{bmatrix} \tau_m^{g,1,1} & \tau_m^{g,1,2} & \tau_m^{g,1,3} \\ \tau_m^{g,2,1} & \tau_m^{g,2,2} & \tau_m^{g,2,3} \\ \tau_m^{g,3,1} & \tau_m^{g,3,2} & \tau_m^{g,3,3} \end{bmatrix}.$$

The elements $\tau_m^{g,r,r'}$ are determined by the contact rate between risk groups r and r' of the opposite sex and the transmission probability per contact $p_{m,0}^{g,r,r'}$ between these risk groups for an infected person in stage m defined in Eq. S17:

$$\tau_m^{g,r,r'}(t) = c^{g,r}(t) \cdot \left(\varepsilon \delta_{r,r'} + (1 - \varepsilon) \frac{c^{g',r'}(t) \pi^{g',r'}}{\sum_r c^{g',r}(t) \pi^{g',r'}} \right) \cdot p_{m,0}^{g,r,r'}. \tag{S34}$$

The matrix $\mathbf{\Sigma}$ includes progression between disease stages (σ_m), movement from higher to lower sexual risk groups (ψ), and removal from the sexually active population ($\nu = 1/35$ per year). As these parameters are assumed fixed and the same for each sex, $\mathbf{\Sigma}$ does not depend on time and consists of identical 15×15 submatrices on the block diagonal

$$\mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}^G & 0 \\ 0 & \mathbf{\Sigma}^G \end{bmatrix}. \tag{S35}$$

The matrix $\mathbf{\Sigma}^G$ consists of 3×3 submatrices $\mathbf{\Sigma}_m^G$ on the diagonal representing transitions between risks groups and removals from each disease stage, and scaled identity matrices on the subdiagonal for entrants into the next disease stage:

$$\mathbf{\Sigma}^G = \begin{bmatrix} \mathbf{\Sigma}_1^G & 0 & 0 & 0 & 0 \\ \sigma_1 \mathbf{I}_3 & \mathbf{\Sigma}_2^G & 0 & 0 & 0 \\ 0 & \sigma_2 \mathbf{I}_3 & \mathbf{\Sigma}_3^G & 0 & 0 \\ 0 & 0 & \sigma_3 \mathbf{I}_3 & \mathbf{\Sigma}_4^G & 0 \\ 0 & 0 & 0 & \sigma_4 \mathbf{I}_3 & \mathbf{\Sigma}_5^G \end{bmatrix}, \tag{S36}$$

where \mathbf{I}_3 is the 3×3 identity matrix and $\mathbf{\Sigma}_m^G = -(\nu + \sigma_m) \mathbf{I}_3 + \mathbf{\Psi}$ with $\mathbf{\Psi} = \begin{bmatrix} -2\psi & 0 & 0 \\ \psi & -\psi & 0 \\ \psi & \psi & 0 \end{bmatrix}$ representing transitions from higher to lower sexual risk groups.

We solve for $R_0(t)$ numerically as the dominant eigenvalue of $\mathbf{K}_L(t) = -\mathbf{T}(t)\mathbf{\Sigma}^{-1}$. Note that because $c^{g,r}(t)$ is scaled proportionally for all (g,r) , it follows from Eq. S1 that

$$\mathbf{T}(t) = \left(1 - \Delta_{\bar{c}} \left(1 - \frac{1}{1 + \exp\left(\frac{t - (t_{\bar{c}} + d_{\bar{c}}/2)}{d_{\bar{c}}/10}\right)} \right) \right) \mathbf{T}(0). \tag{S37}$$

Thus,

$$R_0(t) = \left(1 - \Delta_{\bar{c}} \left(1 - \frac{1}{1 + \exp\left(\frac{t - (t_{\bar{c}} + d_{\bar{c}}/2)}{d_{\bar{c}}/10}\right)} \right) \right) R_0(0), \tag{S38}$$

and in particular, after simulated behavior change has completed at the start of the intervention period $R_0(t)_{t > 2010} \approx (1 - \Delta_{\bar{c}}) R_0(0)$.

Note that following Diekmann et al., we can obtain the NGM \mathbf{K} from the NGM with large domain \mathbf{K}_L by $\mathbf{K} = \mathbf{E}'\mathbf{K}_L\mathbf{E}$, where

$$\mathbf{E}' = \begin{bmatrix} \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \\ \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{I}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 & \mathbf{0}_3 \end{bmatrix}.$$

1. Statistics South Africa (2010) Mid-year population estimates 2010. Available at www.statssa.gov.za/publications/P0302/P03022010.pdf. Accessed March 26, 2012.
2. Johnson LF, Hallett TB, Rehle TM, Dorrington RE (2012) The effect of changes in condom usage and antiretroviral treatment coverage on human immunodeficiency virus incidence in South Africa: A model-based analysis. *J R Soc Interface* 9(72):1544–1554.
3. Hallett TB, Gregson S, Mugurungi O, Gonese E, Garnett GP (2009) Assessing evidence for behaviour change affecting the course of HIV epidemics: A new mathematical modelling approach and application to data from Zimbabwe. *Epidemics* 1(2):108–117.
4. Garnett GP, Anderson RM (1993) Factors controlling the spread of HIV in heterosexual communities in developing countries: Patterns of mixing between different age and sexual activity classes. *Philos Trans R Soc Lond B Biol Sci* 342(1300):137–159.
5. Hollingsworth TD, Anderson RM, Fraser C (2008) HIV-1 transmission, by stage of infection. *J Infect Dis* 198(5):687–693.
6. Wandel S, et al.; eligibility for ART in lower income countries (eART-linc) collaboration (2008) Duration from seroconversion to eligibility for antiretroviral therapy and from ART eligibility to death in adult HIV-infected patients from low and middle-income countries: Collaborative analysis of prospective studies. *Sex Transm Infect* 84(Suppl 1):i31–i36.
7. Wawer MJ, et al. (2005) Rates of HIV-1 transmission per coital act, by stage of HIV-1 infection, in Rakai, Uganda. *J Infect Dis* 191(9):1403–1409.
8. Donnell D, et al.; Partners in Prevention HSV/HIV Transmission Study Team (2010) Heterosexual HIV-1 transmission after initiation of antiretroviral therapy: A prospective cohort analysis. *Lancet* 375(9731):2092–2098.
9. Matthews GV, et al. (2002) Virological suppression at 6 months is related to choice of initial regimen in antiretroviral-naïve patients: A cohort study. *AIDS* 16(1):53–61.
10. Egger M, et al.; ART Cohort Collaboration (2002) Prognosis of HIV-1-infected patients starting highly active antiretroviral therapy: A collaborative analysis of prospective studies. *Lancet* 360(9327):119–129.
11. May M, et al.; IeDEA Southern Africa and West Africa (2010) Prognosis of patients with HIV-1 infection starting antiretroviral therapy in sub-Saharan Africa: A collaborative analysis of scale-up programmes. *Lancet* 376(9739):449–457.
12. Cornell M, et al.; International Epidemiologic Databases to Evaluate AIDS Southern Africa (IeDEA-SA) Collaboration (2010) Temporal changes in programme outcomes among adult patients initiating antiretroviral therapy across South Africa, 2002–2007. *AIDS* 24(14):2263–2270.
13. Van Cutsem G, et al. (2011) Correcting for mortality among patients lost to follow up on antiretroviral therapy in South Africa: A cohort analysis. *PLoS One* 6(2):e14684.
14. Mutevedzi PC, Lessells RJ, Newell ML (2013) Disengagement from care in a decentralised primary health care antiretroviral treatment programme: Cohort study in rural South Africa. *Trop Med Int Health* 18(8):934–941.
15. Fatti G, Grimwood A, Bock P (2010) Better antiretroviral therapy outcomes at primary healthcare facilities: An evaluation of three tiers of ART services in four South African provinces. *PLoS One* 5(9):e12888.
16. Nglazi MD, et al. (2011) Changes in programmatic outcomes during 7 years of scale-up at a community-based antiretroviral treatment service in South Africa. *J Acquir Immune Defic Syndr* 56(1):e1–e8.
17. Boule A, et al. (2010) Seven-year experience of a primary care antiretroviral treatment programme in Khayelitsha, South Africa. *AIDS* 24(4):563–572.
18. Fox MP, et al. (2012) Treatment outcomes after seven years of public-sector HIV treatment at the Themba Lethu Clinic in Johannesburg, South Africa. *AIDS* 26(14):1823–1828.
19. South Africa Department of Health (2011) *The 2010 National Antenatal Sentinel HIV and Syphilis Prevalence Survey in South Africa* (National Department of Health, Pretoria, South Africa).
20. Human Sciences Research Council (2002) *South African National HIV Prevalence, Behavioural Risks and Mass Media Household Survey 2002* (HSRC Press, Cape Town, South Africa).
21. Shisana O, et al. (2005) *South African National HIV Prevalence, HIV Incidence, Behaviour and Communication Survey, 2005* (HSRC Press, Cape Town, South Africa).
22. Shisana O, et al. (2009) *South African National HIV Prevalence, Incidence, Behaviour and Communication Survey 2008: A Turning Tide Amongst Teenagers?* (HSRC Press, Cape Town, South Africa).
23. Department of Health Republic of South Africa (2011) *Department of Health Republic of South Africa. National Strategic Plan for HIV and AIDS/CCMT Monthly Statistics June 2011* (National Department of Health, Pretoria, South Africa).
24. Johnson L, Dorrington R, Bradshaw D, Pillay-Van Wyk V, Rehle T (2009) Sexual behaviour patterns in South Africa and their association with the spread of HIV: Insights from a mathematical model. *Demogr Res* 21:289–340.
25. Raftery AE, Bao L (2010) Estimating and projecting trends in HIV/AIDS generalized epidemics using incremental mixture importance sampling. *Biometrics* 66(4):1162–1173.
26. Diekmann O, Heesterbeek JAP, Roberts MG (2010) The construction of next-generation matrices for compartmental epidemic models. *J R Soc Interface* 7(47):873–885.

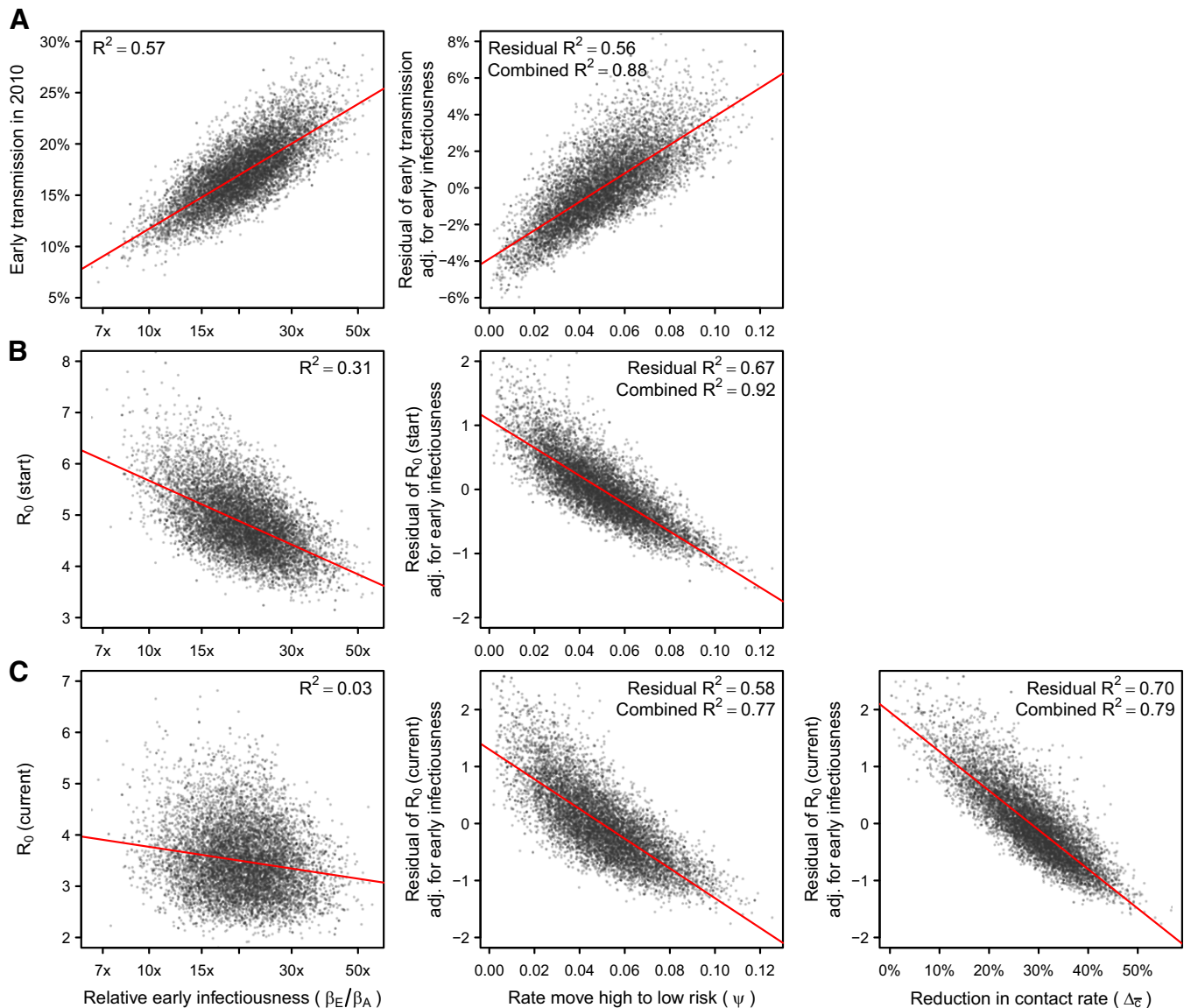


Fig. S2. The relationship between (A) early transmission in 2010, (B) R_0 at the start of the epidemic, and (C) R_0 during the intervention period and model parameters: log of relative infectiousness during early infection (Left; β_E/β_A in Table S5), the rate of movement from higher- to lower-risk groups (Center; ψ in Table S5), and the percentage reduction in sexual contact rate over time (Bottom right; $\Delta_{\bar{c}}$ in Table S5). The leftmost plot indicates the relationship between the outcome and relative early infectiousness (β_E/β_A). Center and Right indicate the relationship between the residual variation in the outcome (after removing the effect of β_E/β_A) and ψ (Center) or $\Delta_{\bar{c}}$ (Right). Residual R^2 values indicate the fraction of residual variance explained by ψ (Center) and $\Delta_{\bar{c}}$ (Right), and combined R^2 values indicate the fraction of all variance explained by early infectiousness and the plotted parameter together. All three parameters together explain 95% of the variation in R_0 during the intervention period (C).

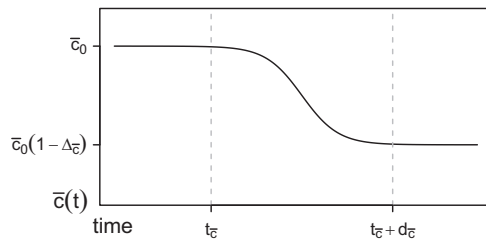


Fig. S7. Logistic curve describing behavior change in Eq. S1.

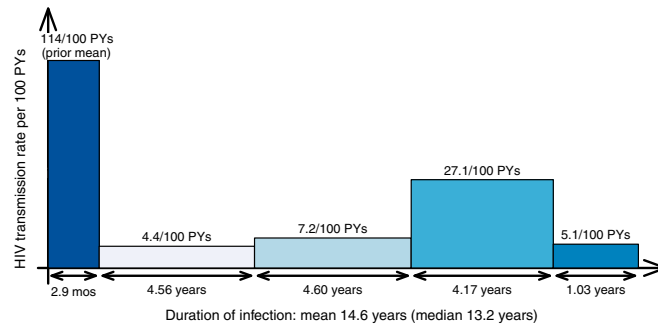


Fig. S8. Average duration of and relative HIV transmission rate during each stage of HIV infection. (Note that vertical axis is not drawn to scale.)

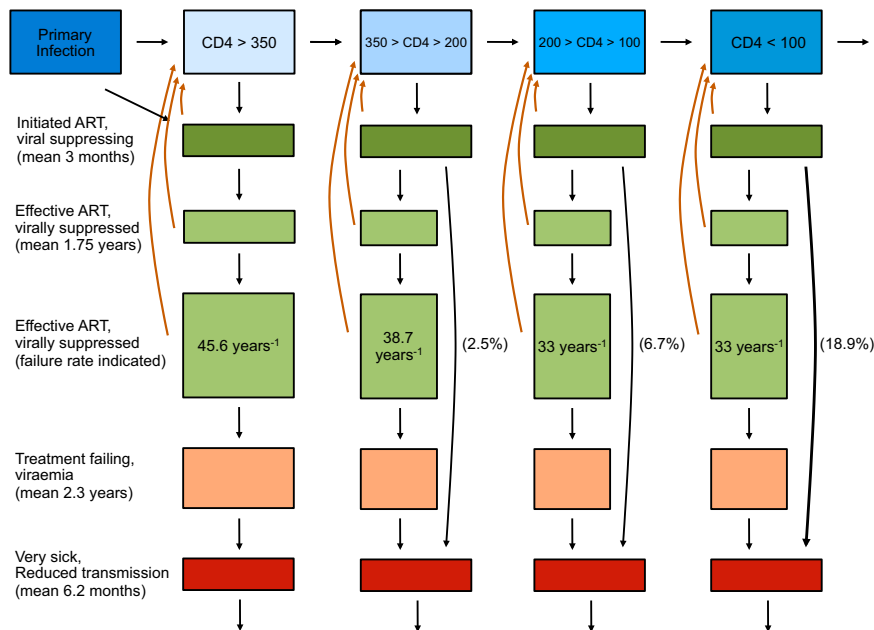


Fig. S9. Stages of ART. Arrows indicate possible movements of individuals between stages. Black arrows indicate natural progression of individuals. Red arrows indicate dropout from ART. Individuals may drop out from ART only once and may return to a different CD4 stage depending on their duration on ART (SI Appendix, Tables S4 and S5).

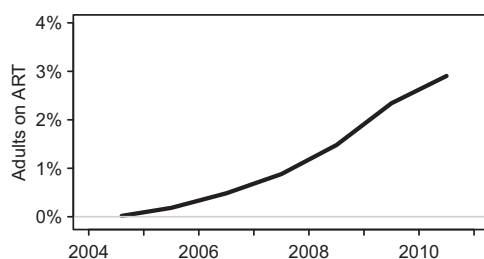


Fig. S10. The percentage of adults (aged 15+) on ART over time used in the model calibration, based on reported numbers on ART from the South Africa Department of Health and population size estimates from Statistics South Africa.

Table S1. Stages of ART

Subscript	Stage	Duration	Infectiousness
0	Untreated, ART naïve		
1	Virally suppressing	3 mo	50% lower than prev. stage
2	Early effective ART, virally suppressed	1.75 y	92% lower than CD4 ≤ 350
3	Effective ART, virally suppressed	(Fig. S9)	92% lower than CD4 ≤ 350
4	Treatment failing, viremic	2.3 y	Same as CD4 200–350
5	Very sick	6.2 mo	Same as CD4 ≤ 100
6	Untreated, dropped out after first initiation, eligible to restart ART		
7	Reinitiated ART, virally suppressing	3 mo	50% lower than prev. stage
8	Reinitiated ART, early effective ART	1.75 y	92% lower than CD4 ≤ 350
9	Reinitiated ART, effective ART	(Fig. S9)	92% lower than CD4 ≤ 350
10	Reinitiated ART, treatment failing, viremic	2.3 y	Same as CD4 100–200
11	Reinitiated ART, very sick	6.2 mo	Same as CD4 ≤ 100
12	Untreated, dropped out after second initiation, not eligible to restart		

Table S2. Rate per year of dropping out from ART

Baseline CD4 cell count	Duration on treatment			
	Virally suppressing	Early effective ART	Effective ART	Reinitiated ART
>350	0.168	0.168	0.088	0.06
200–350	0.168	0.168	0.088	0.06
100–200	0.156	0.156	0.088	0.06
≤ 100	0.120	0.120	0.088	0.06

Table S3. CD4 stage after dropping out of treatment

CD4 stage at ART initiation	Treatment stage at dropout	Dropouts returning to CD4 category, %			
		CD4 >350	CD4 200–350	CD4 100–200	CD4 ≤100
>350	Virally suppressing	100			
	Early effective ART	100			
	Effective ART	100			
200–350	Virally suppressing		100		
	Early effective ART	100			
	Effective ART	100			
100–200	Virally suppressing			100	
	Early effective ART	50	50		
	Effective ART	100			
≤100	Virally suppressing				100
	Early effective ART		50	50	
	Effective ART		100		

Table S4. Model parameters

Parameter	Description	Value
α	Population growth rate (in absence of HIV)	0.023 per year
ν	Rate of progression from 15–49 to 50+ age groups	1/35 per year
μ	Mortality rate out of the 50+ age group	1/11.45 per year
$\pi^{g,r}$	Proportion of the age 15–49 population of sex g in risk group r in absence of HIV	Estimated
$\bar{\pi}^{g,r}$	Proportion of new entrants of sex g entering risk group r	Derived from π and ψ
$\Psi_{r,r'}^g$	Annual rate of moving from risk group r to r' for sex g	Estimated
$\bar{c}(t)$	Population mean contact rate per year at time t	Estimated
$\omega^{g,r}$	Relative contact rate between risk group r and low-risk group for sex g	Estimated
ε	Degree of assortative mixing	Estimated
θ_G	Balance between male and female partner preference	0.5
κ_{r_M,r_F}	Intensity of partnership between male in risk group r_M and female in r_F	Estimated
$\beta_{m,u}$	Annual HIV transmission rate in stage m and ART status u	See Figs. S8 and S9
σ_m	Rate of progression from HIV stage m to stage $m+1$	$\begin{bmatrix} 0.24 \\ 4.56 \\ 4.53 \\ 4.28 \\ 0.94 \end{bmatrix} \text{ year}^{-1}$
$\bar{\sigma}_m$	Rate of progression from HIV stage m to stage $m+1$ after treatment dropout	2σ
λ_m^g	Rate of ART initiation for sex g in stage m	See text
$\tilde{\lambda}_m^g$	Rate of reinitiating ART after dropout for sex g in stage m	SI Appendix, section 1.5
$\phi_{m,u}^g$	Rate of progression from ART stage u to $u+1$ when initiating ART in HIV stage m for sex g	Fig. S9
ξ_m	Probability of immediate treatment failure if initiating ART in stage m	$\begin{bmatrix} 0 \\ 0 \\ 0.025 \\ 0.067 \\ 0.189 \end{bmatrix}$
$\eta^{m,u}$	Rate of dropping out of ART if initiated in stage m and currently in stage u	Table S2
$\bar{\eta}^{m,u}$	Rate of dropping out of ART after reinitiating if reinitiated in stage m and currently in stage u	Table S2
$\rho_{m,u}^{m',u'}$	Probability of entering CD4 stage m after dropping out of stage (m',u')	Table S3
t_0	Date at which HIV epidemic is seeded into population	Estimated
$\delta_m^{g,r}$	Seed HIV prevalence for sex g , risk group r , and disease stage m	SI Appendix, section 1.7

Table S5. Estimated model parameters and prior distributions

Parameter	Description	Prior
t_0	Start date of the epidemic	Unif(1983, 1988)
$1 - \pi^{M,L}$	Proportion of males not in the low-risk group	Unif(0.05, 0.7)
$1 - \pi^{F,L}$	Proportion of females not in the low-risk group	Unif(0.05, 0.7)
$\frac{\pi^{M,H}}{1 - \pi^{M,L}}$	Proportion of males in high-risk group of those not in low-risk group	Unif(0.2, 0.8)
$\frac{\pi^{F,H}}{1 - \pi^{F,L}}$	Proportion of females in high-risk group of those not in low-risk group	Unif(0.2, 0.8)
Ψ	Annual rate of movement from higher- to lower-risk groups	Unif(0.0, 0.15)
\bar{c}_0	Mean annual contact rate at start of the epidemic	Unif(0.5, 4.0)
$\Delta\bar{c}$	Proportion reduction in average contact rate	Unif(0.0, 0.7)
$t_{\bar{c}}$	Year behavior change starts	Unif(1990, 2002)
$t_{\bar{c}} + d_{\bar{c}}$	Year behavior change ends	Unif(2002, 2010)
$\omega^{F,M}$	Relative contact rate between medium-risk and low-risk females	Unif(1, 70)
$\omega^{F,H} - \omega^{F,M}$	Additional relative contact rate for high-risk women	Unif(0, 50)
ε	Assortativity of sexual mixing	Unif(0.2, 0.8)
κ_H	Partnership intensity for partnership involving a high-risk partner	*
κ_M	Partnership intensity for partnership between medium and low risk	*
κ_L	Partnership intensity for partnership between low-risk partners	*
β_E/β_A	Ratio of infectiousness during early and asymptomatic (CD4 >350) infection	LogNormal(3.2, 0.34)
γ	Bias in ANC prevalence and 15–49-y-old female prevalence on the logit scale	Unif(0.0, 1.0)

*These parameters have a joint uniform prior distribution such that $0 < \kappa_L < \kappa_M < \kappa_H < 1$.

Table S6. Bivariate correlations between posterior parameter estimates

Parameter	t_0	Risk group size				ψ	Contact rate				ω	ε	Partnership intensity		Early infectious	
t_0																
$1 - \pi^{M,L}$	-0.0															
$1 - \pi^{F,L}$	0.1	0.7														
$\pi^{M,H}/(1 - \pi^{M,L})$	0.0	0.5	0.1													
$\pi^{F,H}/(1 - \pi^{F,L})$	0.0	-0.3	0.1	0.1												
Ψ	0.3	0.3	0.5	-0.0	0.0											
\bar{c}_0	-0.1	0.1	0.2	-0.0	0.2	-0.2										
$\Delta\bar{c}$	0.2	0.2	0.4	-0.0	0.0	0.7	-0.1									
$t_{\bar{c}}$	0.0	0.0	0.0	0.0	-0.0	0.0	-0.0	0.2								
$t_{\bar{c}} + d_{\bar{c}}$	0.0	-0.0	-0.0	-0.0	0.0	0.0	-0.0	0.3	-0.2							
$\omega^{F,M}$	-0.0	0.1	0.1	-0.0	0.1	0.2	-0.1	0.0	0.0							
$\omega^{F,H} - \omega^{F,M}$	0.0	0.0	-0.0	0.1	-0.1	0.1	0.1	0.0	-0.0	0.0	0.1					
ε	-0.1	0.3	0.3	0.0	-0.1	0.1	0.0	0.0	0.0	-0.1	0.0					
κ_H	0.1	-0.5	-0.5	0.1	0.1	0.1	-0.6	0.1	0.0	0.0	0.1	-0.1	-0.1			
κ_M	-0.1	0.0	0.1	-0.0	0.2	0.1	-0.6	0.1	0.0	0.0	-0.0	0.1	-0.0	0.2		
κ_L	-0.1	-0.0	0.1	-0.0	0.1	0.1	-0.3	0.0	0.0	0.0	-0.0	0.1	-0.0	0.1	0.6	
β_E/β_A	0.1	0.1	0.0	0.0	-0.1	-0.5	-0.1	-0.3	-0.1	-0.1	-0.1	-0.0	-0.0	-0.2	-0.1	-0.1

Parameters defined in Table S5.