



# Flow models to interpret population-based deep-sequence pathogen data

Oliver Ratmann (Imperial College London)

Mélodie Monod\*, Andrea Brizzi\*, Ronald M Galiwango\*, Robert Ssekubugu, Yu Chen, Xiaoyue Xi\*, Edward Nelson Kankaka\*, Victor Ssempijja\*, Lucie Abeler Dörner, Adam Akullian, Stella Alamo, Alexandra Blenkinsop, David Bonsall, Larry W Chang, Shozen Dan, Christophe Fraser, Tanya Golubchik, Ronald J Gray, Mathew Hall, Jade Jackson, Godfrey Kigozi, Oliver Laeyendecker, Lisa Mills, Lisa Nelson, Thomas C Quinn, Steven J Reynolds, John Santelli, Nelson Sewankambo, Simon Spencer, Joseph Ssekasanvu, Laura Thompson, Maria J Wawer, David Serwadda, Peter Godfrey-Fausset\*, Joseph Kagaayi, Kate Grabowski\* on behalf of the PANGEA-HIV consortium

Imperial College  
London



Rakai Health  
Sciences Program  
Improved Health Through Research

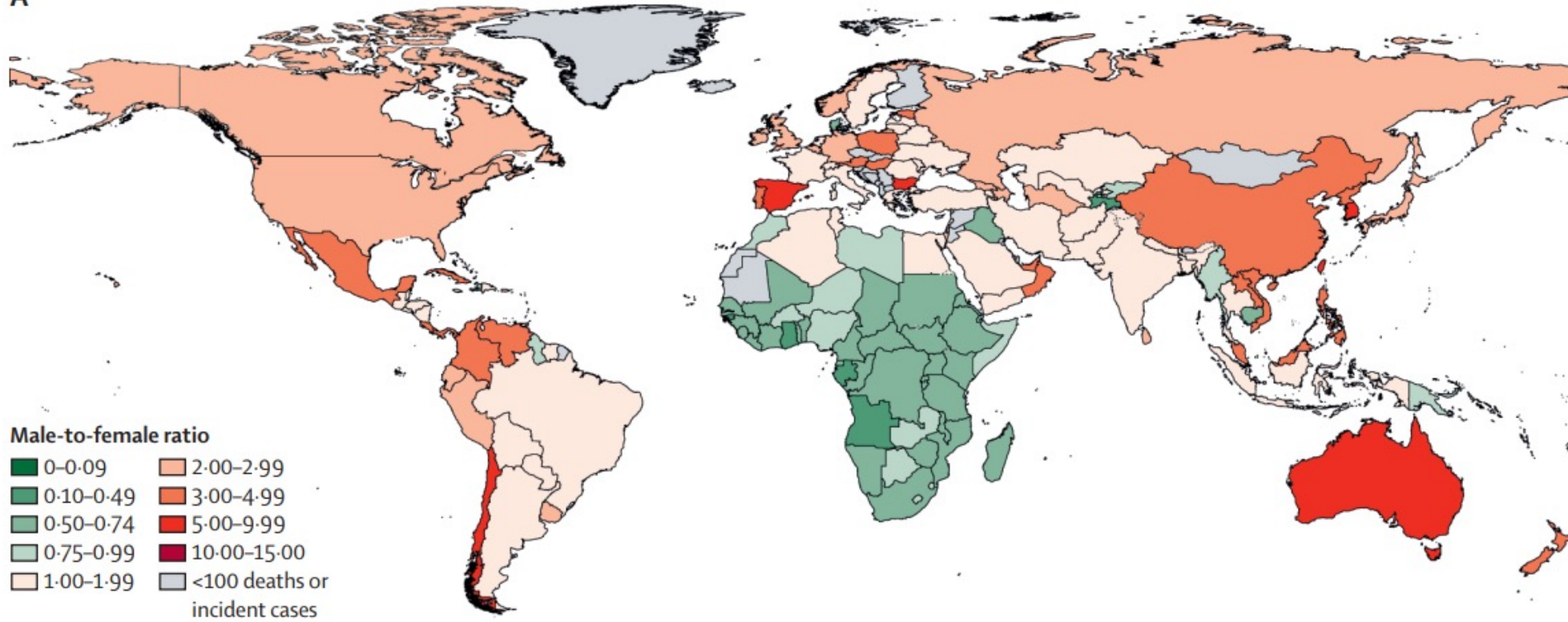
BILL & MELINDA  
GATES foundation



JOHNS HOPKINS  
SCHOOL of MEDICINE

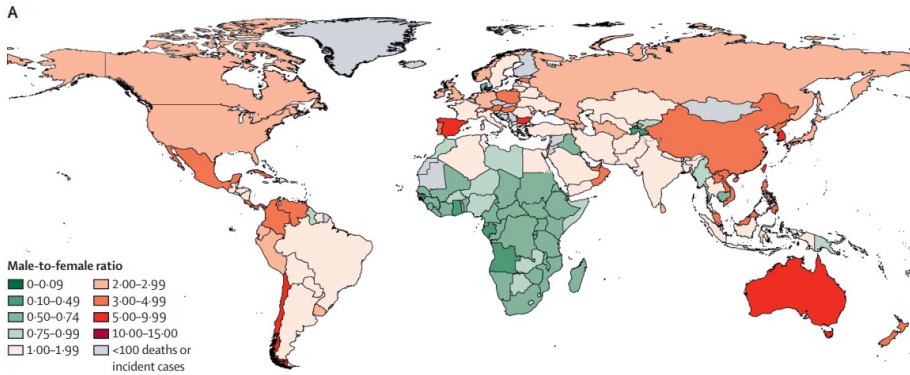


A

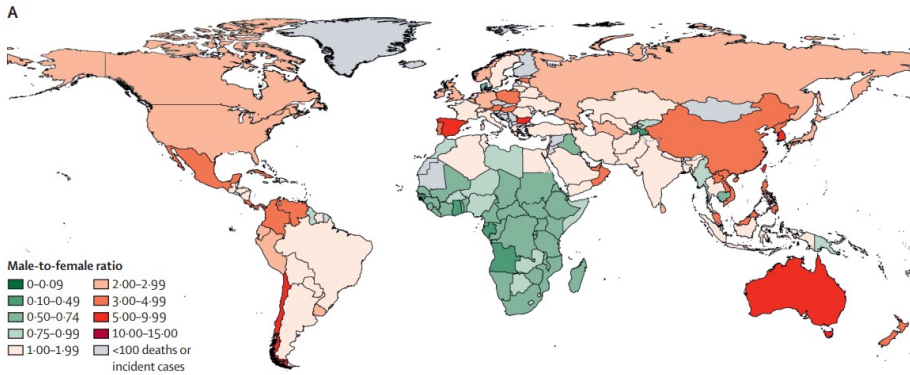


Historically, the African HIV epidemic has been female.





1. What are the recent trends in HIV incidence in women?
2. Are disparities between men and women closing or widening?
3. Which male populations drive incidence in women, and vice versa?
4. What are the best strategies to close gaps and improve population health?

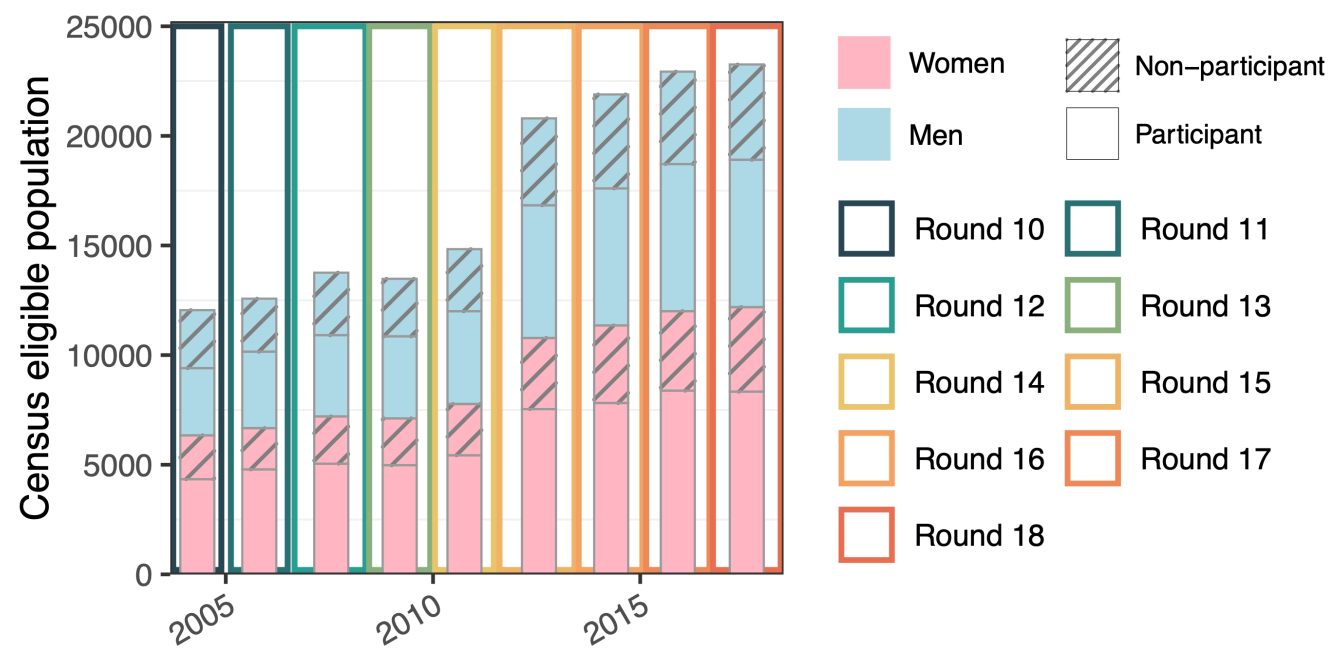
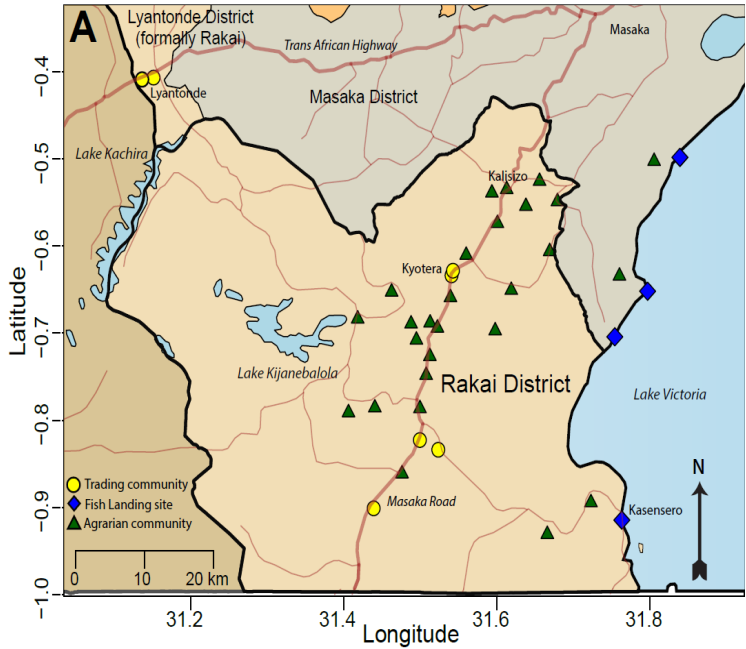


- Data from the Rakai Community Cohort Study
- Longitudinal surveillance of HIV incidence and transmission sources, 2008-2018
- Population-based cohort
- Lower-risk inland communities and high-risk fishing communities; here focus on inland

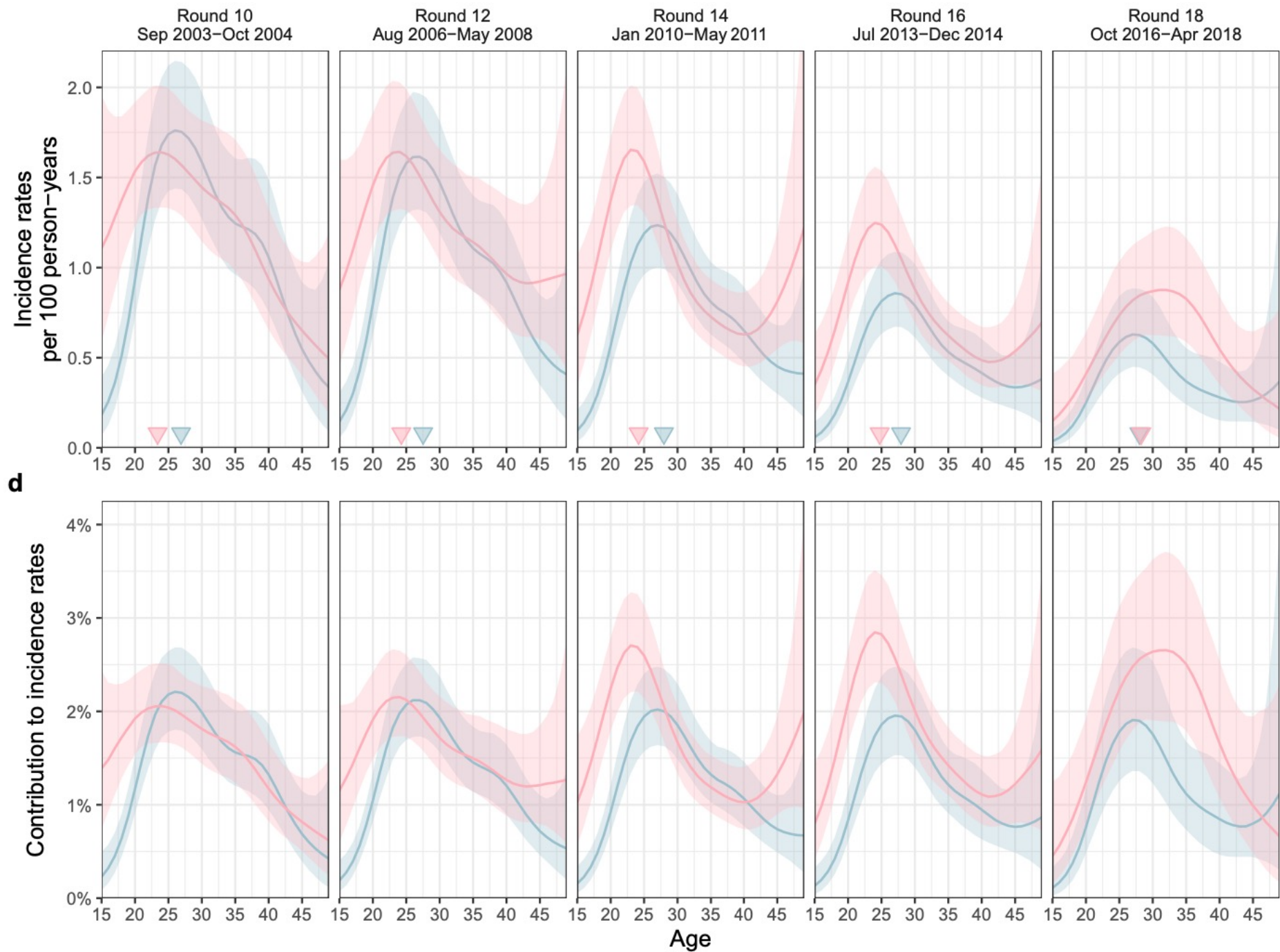


Trends in HIV incidence, 2010 - 2018

• **Rakai  
Community  
Cohort Study**

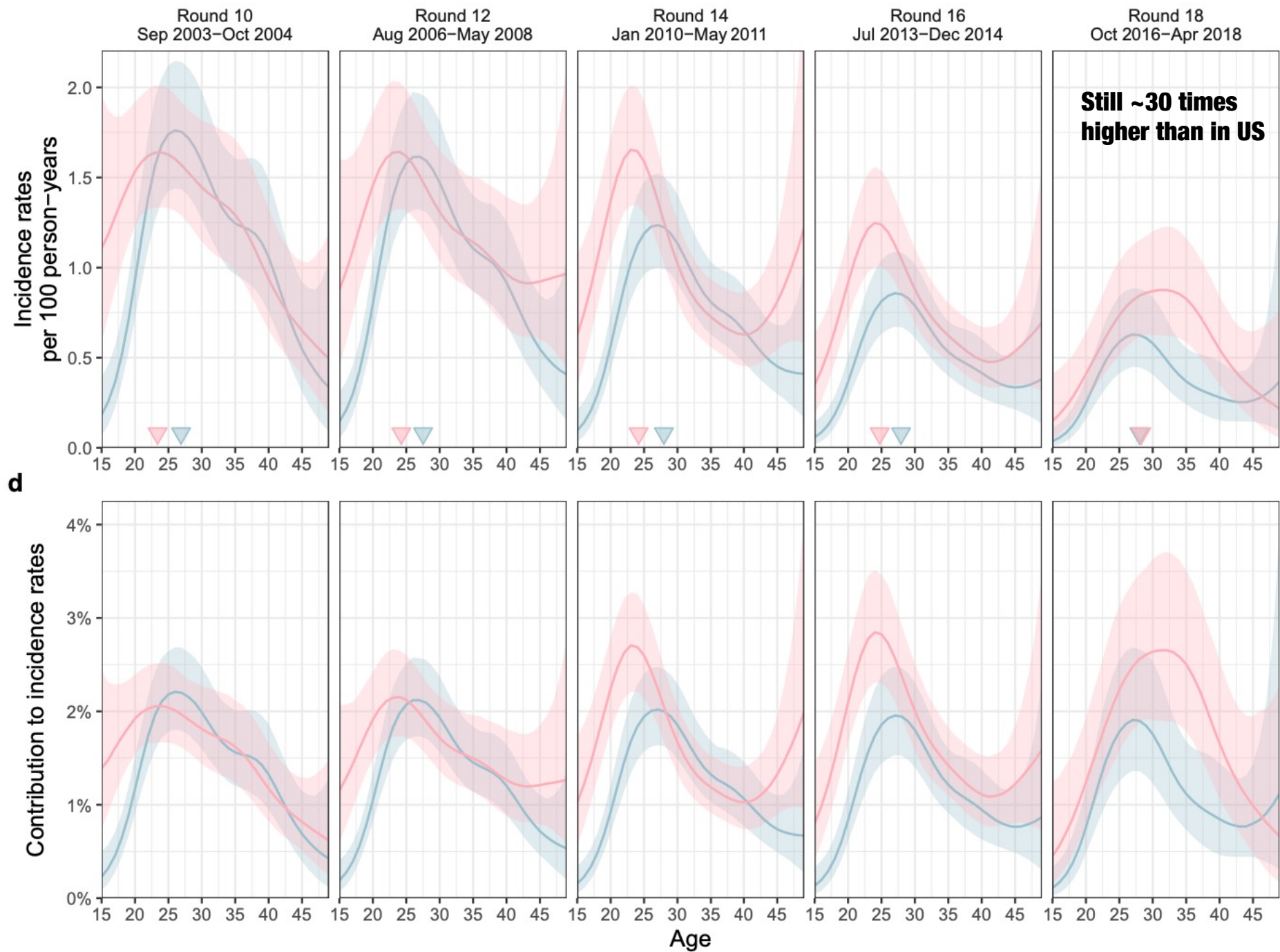


- **1100 incident cases observed over 127k PY, 2003-2018**
- **Faster declines in HIV incidence in men than women, ages 25 and above.**

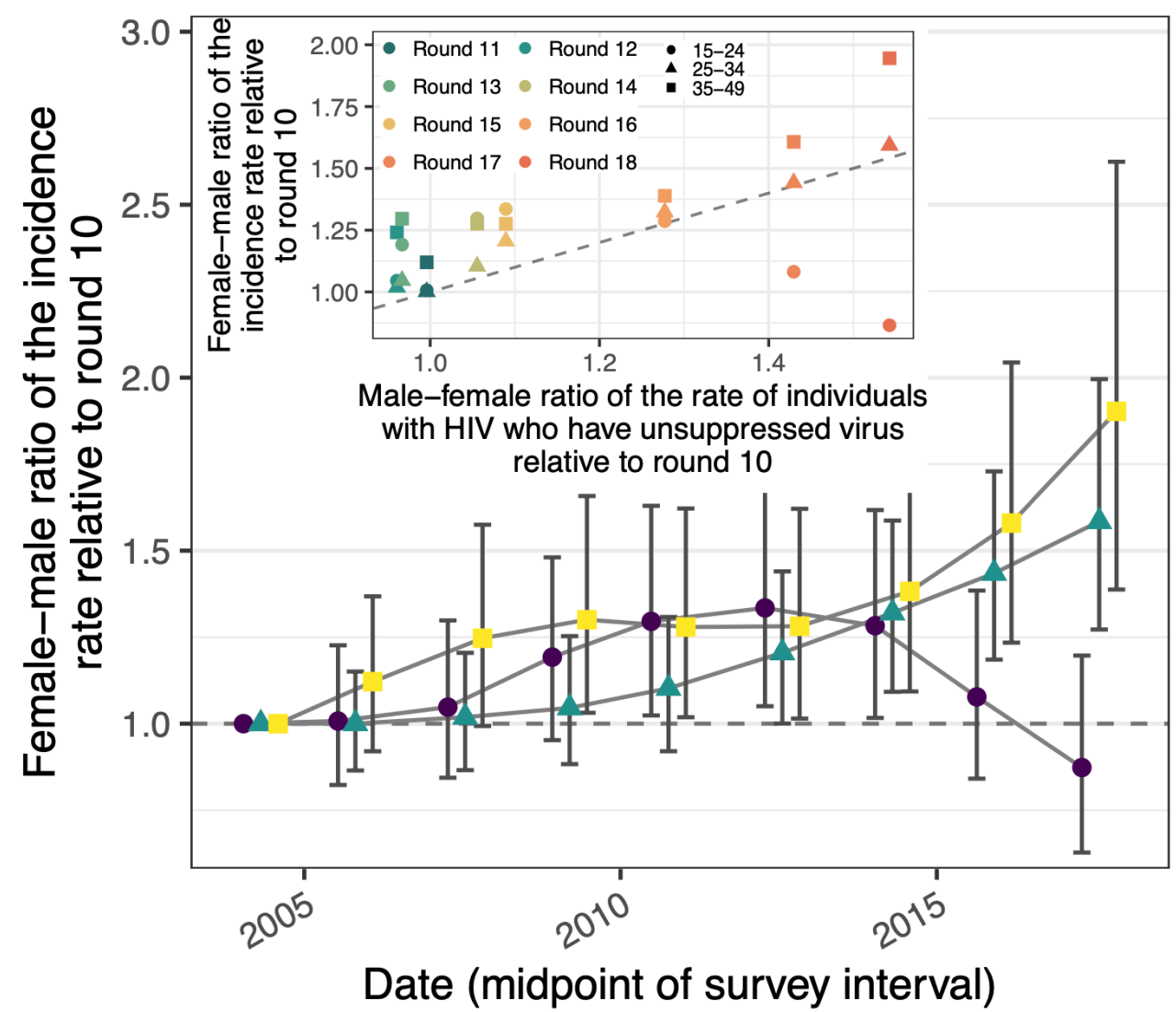
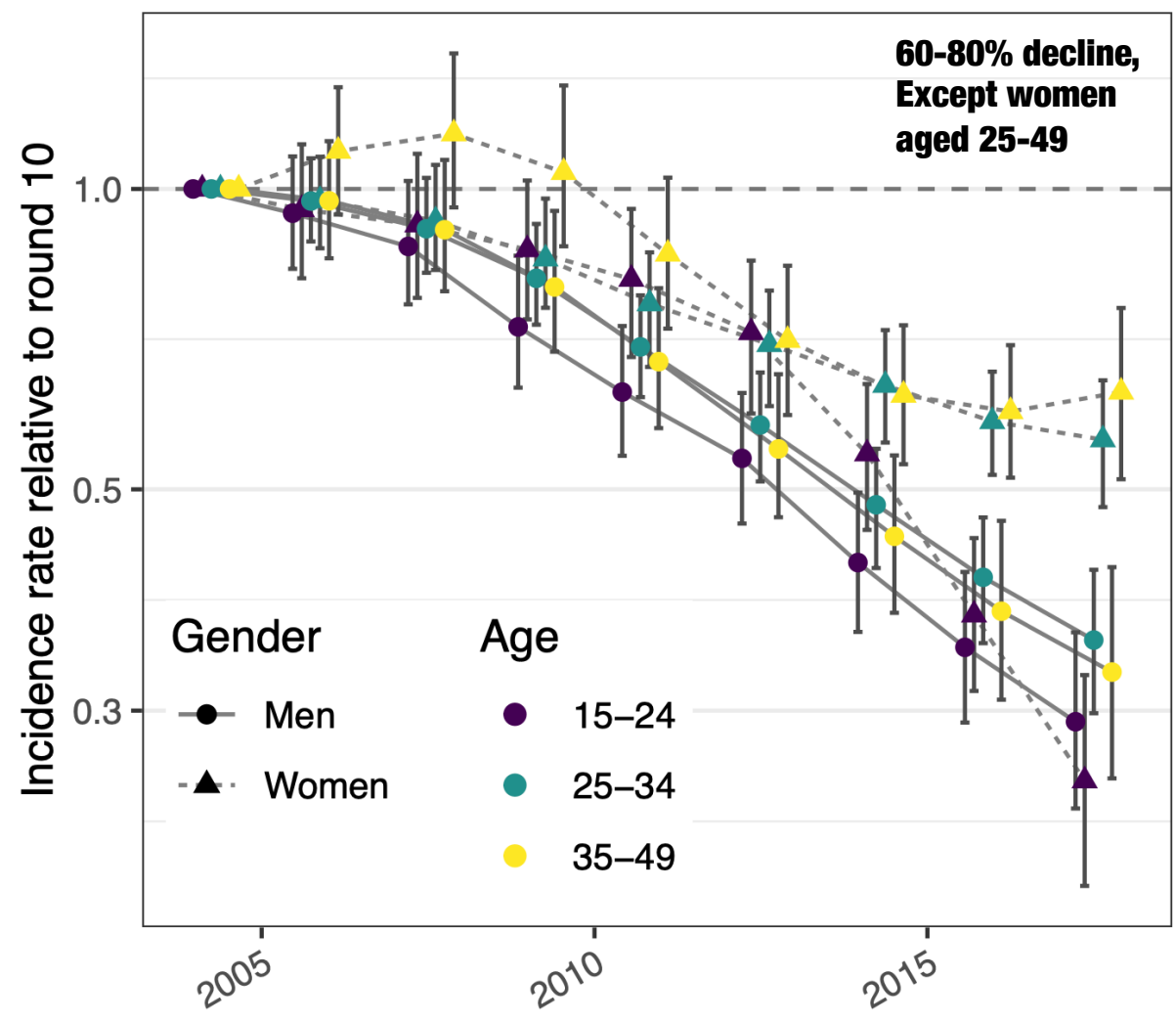




- **1100 incident cases observed over 127k PY, 2003-2018**
- **Faster declines in HIV incidence in men than women, ages 25 and above.**



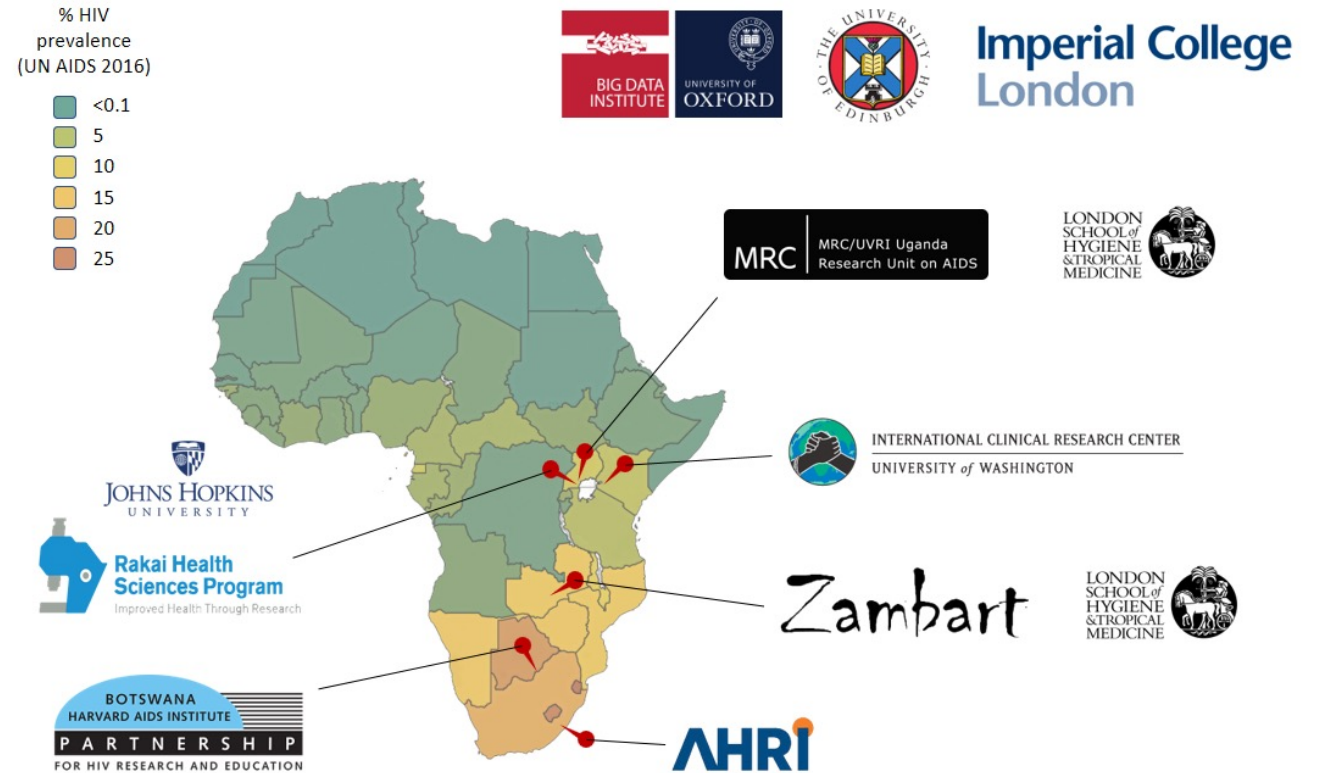
b



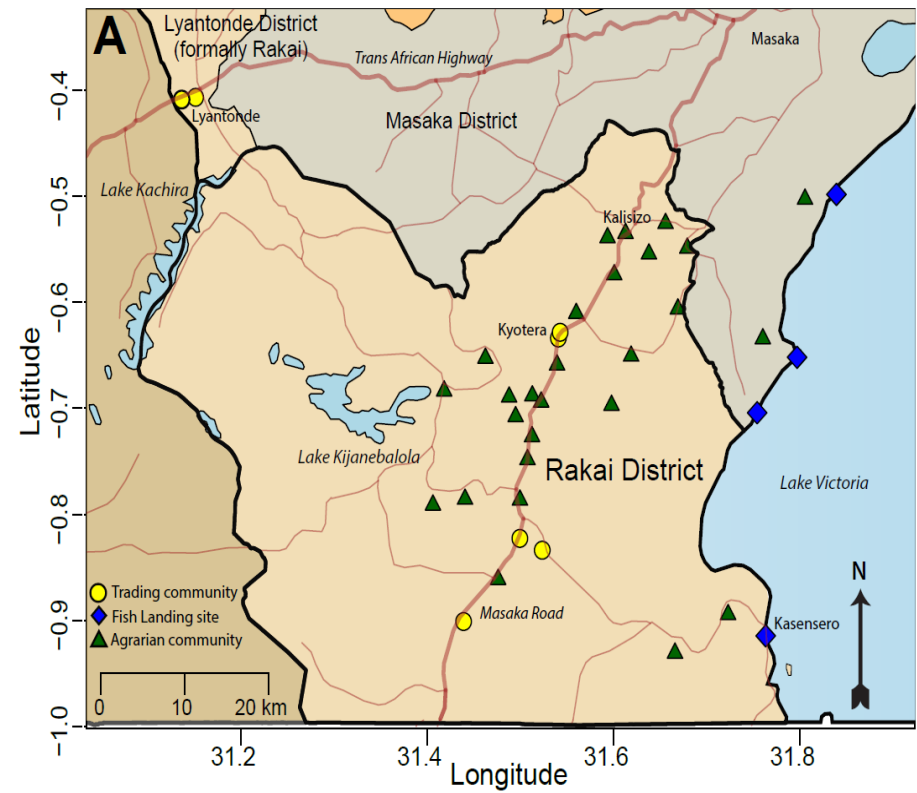
Understanding the changing sources  
of these infections,  
2010 - 2018



- **PANGAEA-HIV:**  
pan-African HIV pathogen  
genomics program  
integrated with population  
surveillance and clinical  
care

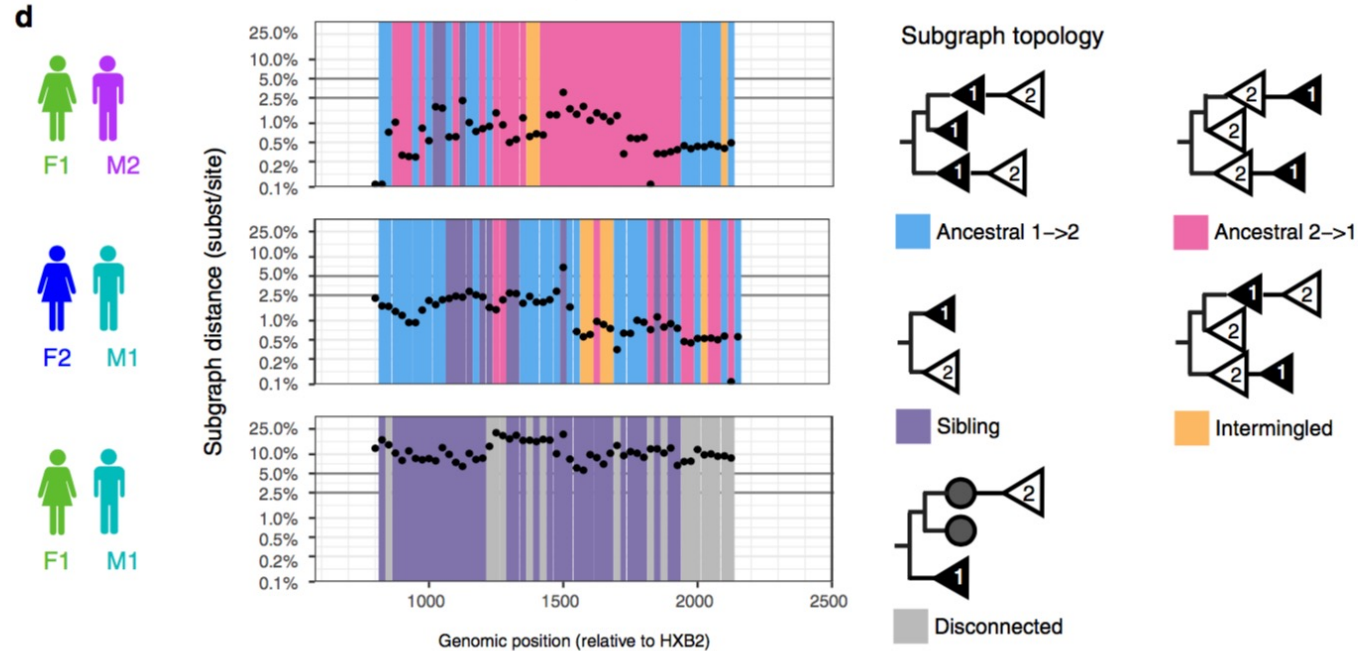
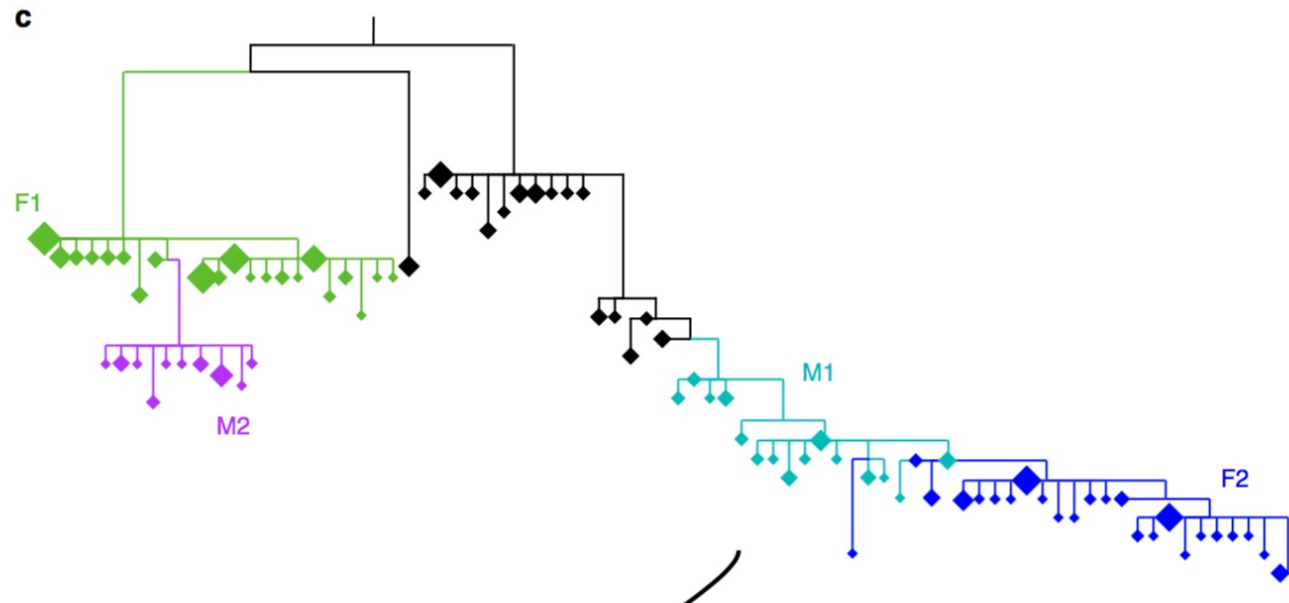


- PANGEA-HIV:**  
pan-African HIV pathogen  
genomics program  
integrated with population  
surveillance and clinical  
care



	Participants with HIV	Participants with HIV reporting no ART use at first visit	Participants with HIV and with virus ever deep-sequenced <sup>†</sup>	
	(n)	(n)	(n)	(%)
Total	5682	4341	2174	38 %
Female (Total)	3817	2836	1291	34 %
Age				
15-24	1066	817	424	40 %
25-34	2074	1488	740	36 %
35-49	1446	826	411	28 %
Male (Total)	1865	1506	883	47 %
Age				
15-24	272	220	157	58 %
25-34	955	782	499	52 %
35-49	984	670	436	44 %
Round <sup>‡</sup>				
10	884	—	115	13 %
11	1002	884	176	18 %
12	1105	912	234	21 %
13	1160	900	368	32 %
14	1741	1392	820	47 %
15	1944	1331	1085	56 %
16	1875	868	892	48 %
17	2015	646	933	46 %
18	1860	432	848	46 %

<sup>†</sup> Individuals with virus ever deep-sequenced were defined as HIV-positive individuals with deep-sequence output meeting minimum quality criteria, see Methods. <sup>‡</sup> Totals by round include individuals seen in other rounds.



## Key discovery of PANGEA-HIV

- HIV deep sequencing provides multiple sequence fragments per person
- Think: phylogeography between individuals
- Inference of transmission direction

Wymant et al. MBE 2017

Hall et al. Elife 2019

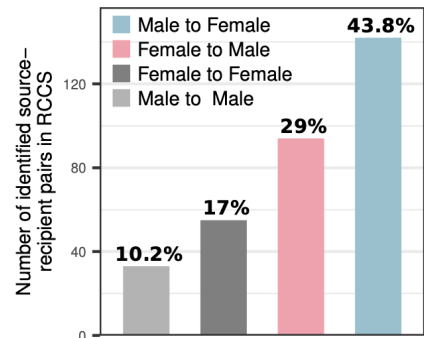
Ratmann et al. Nature Communications 2019

Ratmann et al. Lancet HIV 2020

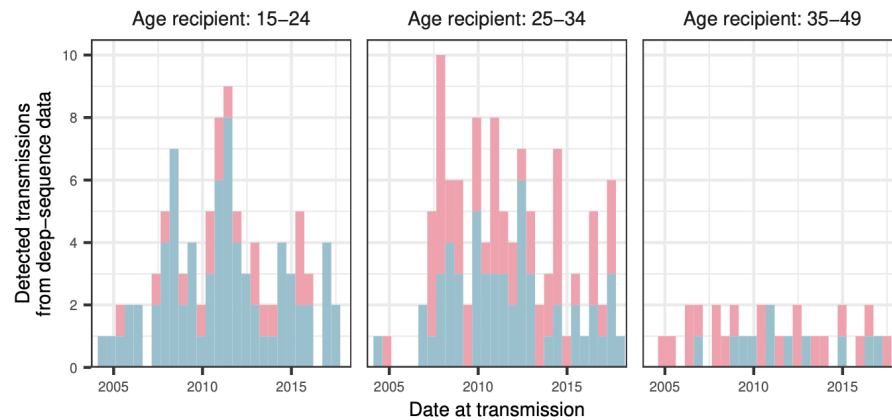
Xi et al. JRSSC 2022



a

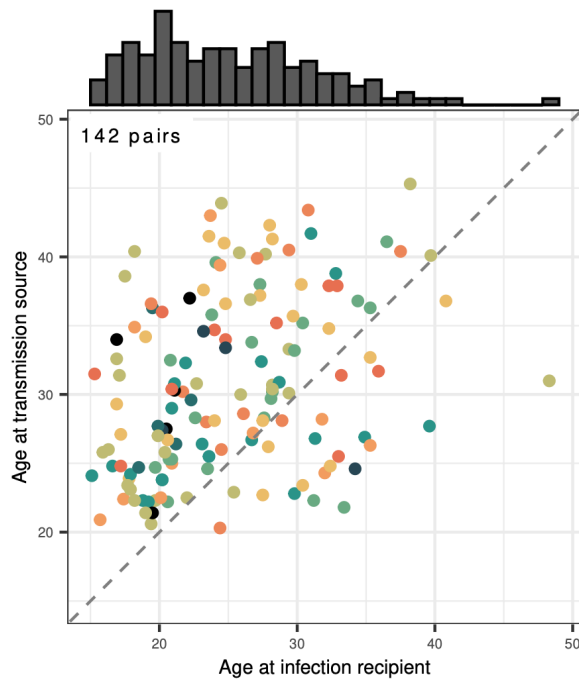


b

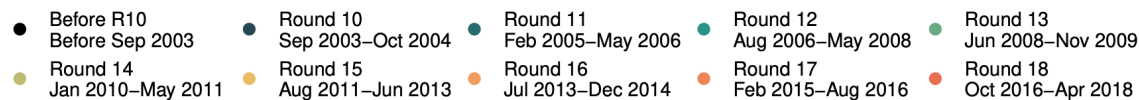
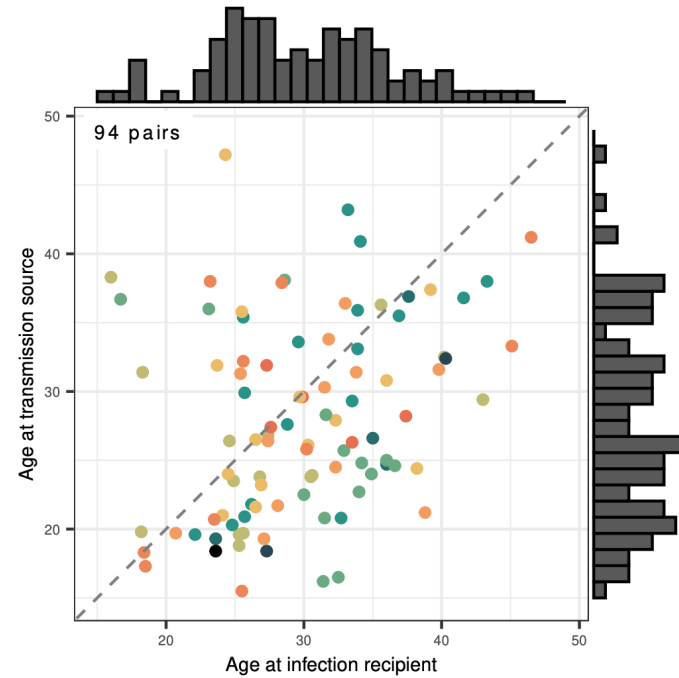


c

Male to female



Female to male



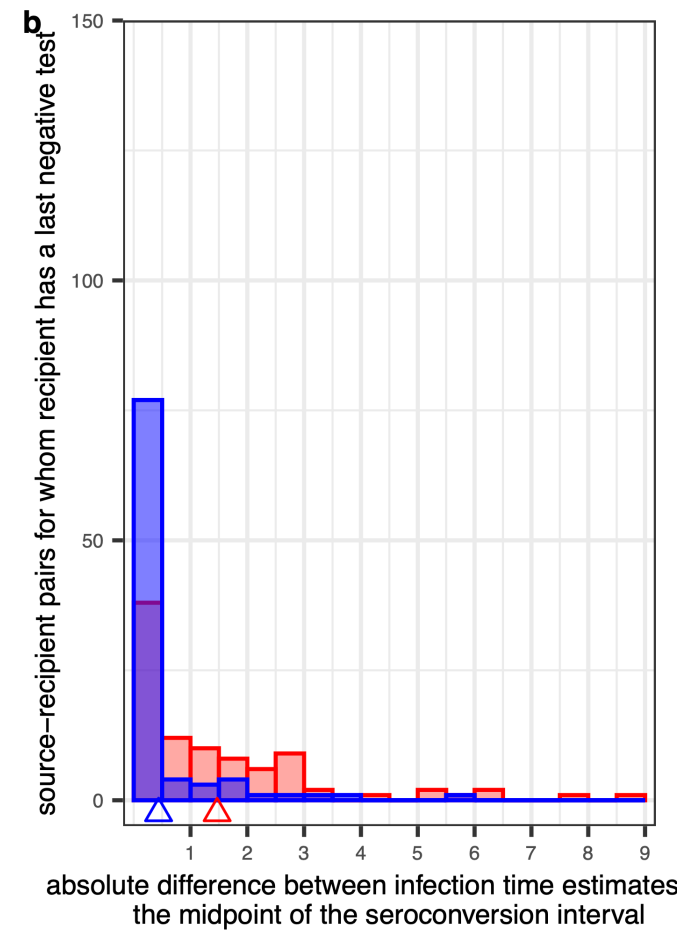
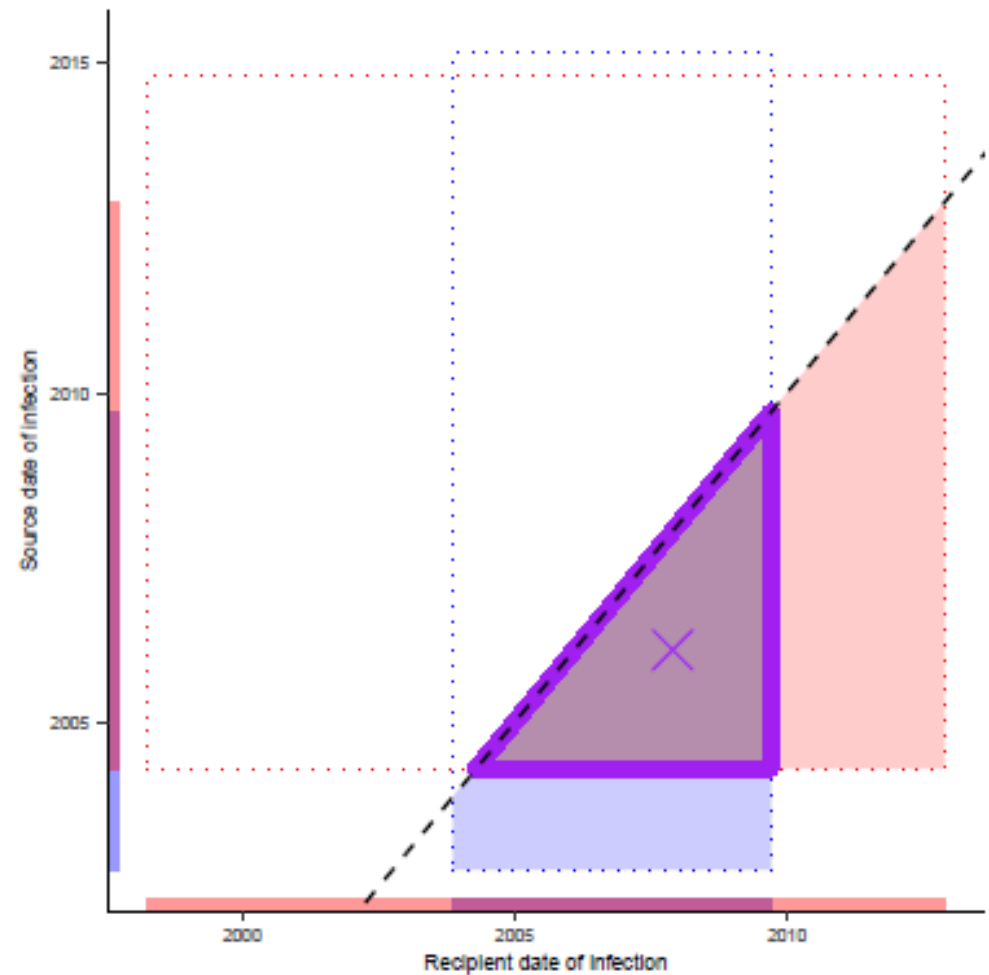
## Transmission cohort, 2013–2018

Identified 236 heterosexual source-recipient pairs

Retained 227 in whom transmission was estimated to have occurred during the study period.

# Dating the likely infection time with deep-sequence data

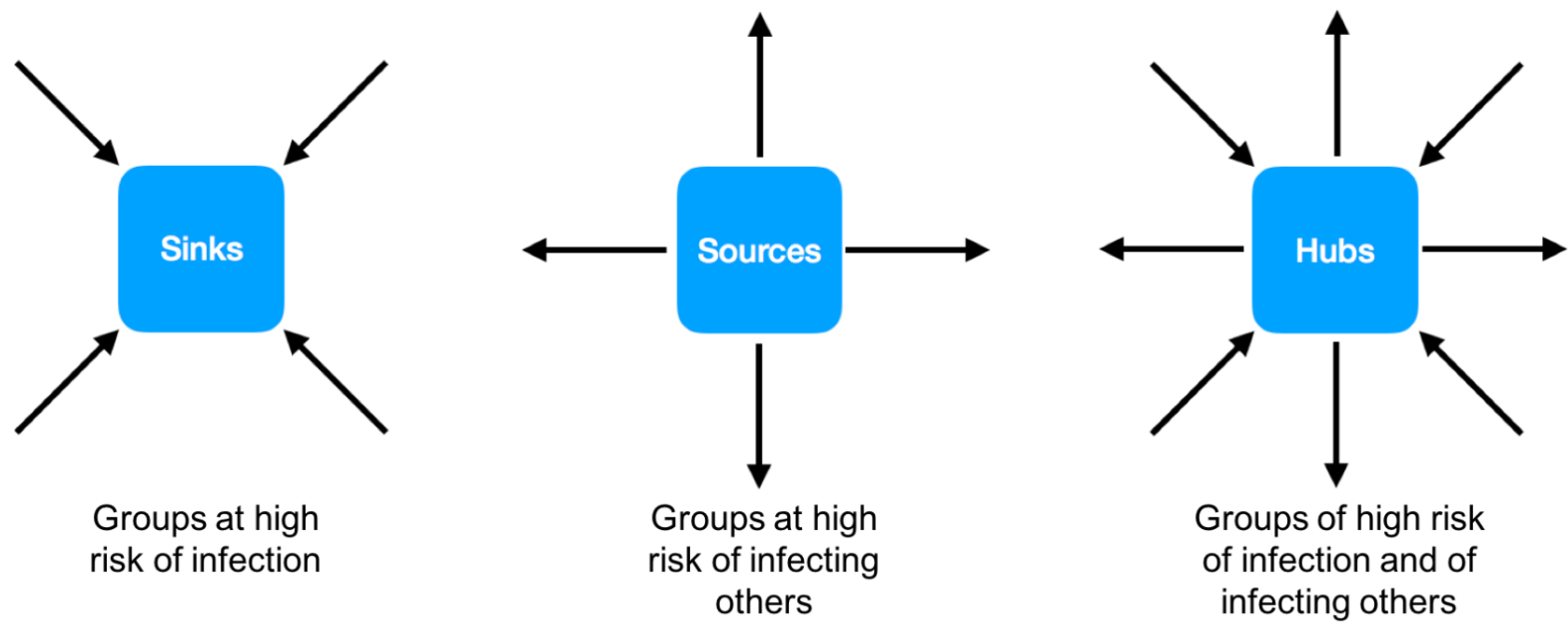
- Used phyloTSI algorithm
- Augmented infection time estimates with epidemiologic data



A statistical perspective on phylodynamics:  
regression models on observed flows



# Concept: transmission sinks, sources, hubs



# Target quantities

- Directional transmission flows,  
e.g. from and to areas with high (h) and low (l) HIV prevalence

$$\boldsymbol{\pi} = \begin{pmatrix} \pi_{hh} & \pi_{hl} \\ \pi_{lh} & \pi_{ll} \end{pmatrix}, \quad \sum_{a,b} \pi_{a,b} = 1$$

# Target quantities

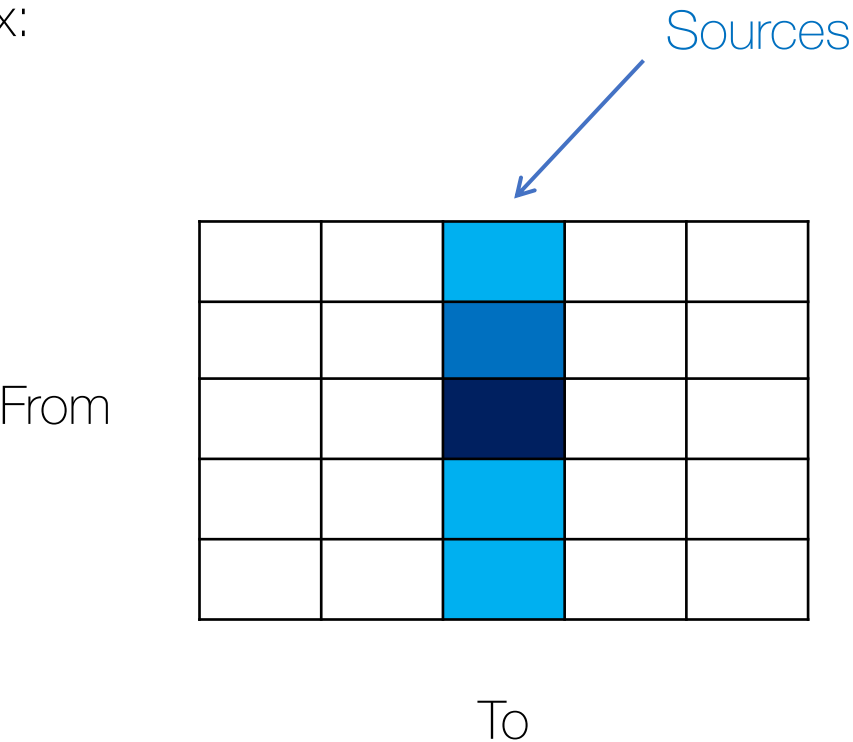
- Directional transmission flows,  
e.g. from and to different age bands:

$$\boldsymbol{\pi} = \begin{pmatrix} \boldsymbol{\pi}^{mf} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\pi}^{fm} \end{pmatrix}, \quad \boldsymbol{\pi}^{mf} = \begin{pmatrix} \pi_{11}^{mf} & \cdots & \pi_{1K}^{mf} \\ \vdots & \ddots & \vdots \\ \pi_{K1}^{mf} & \cdots & \pi_{KK}^{mf} \end{pmatrix} \quad \boldsymbol{\pi}^{fm} = \begin{pmatrix} \pi_{11}^{fm} & \cdots & \pi_{1K}^{fm} \\ \vdots & \ddots & \vdots \\ \pi_{K1}^{fm} & \cdots & \pi_{KK}^{fm} \end{pmatrix}$$

male→female                      female→male

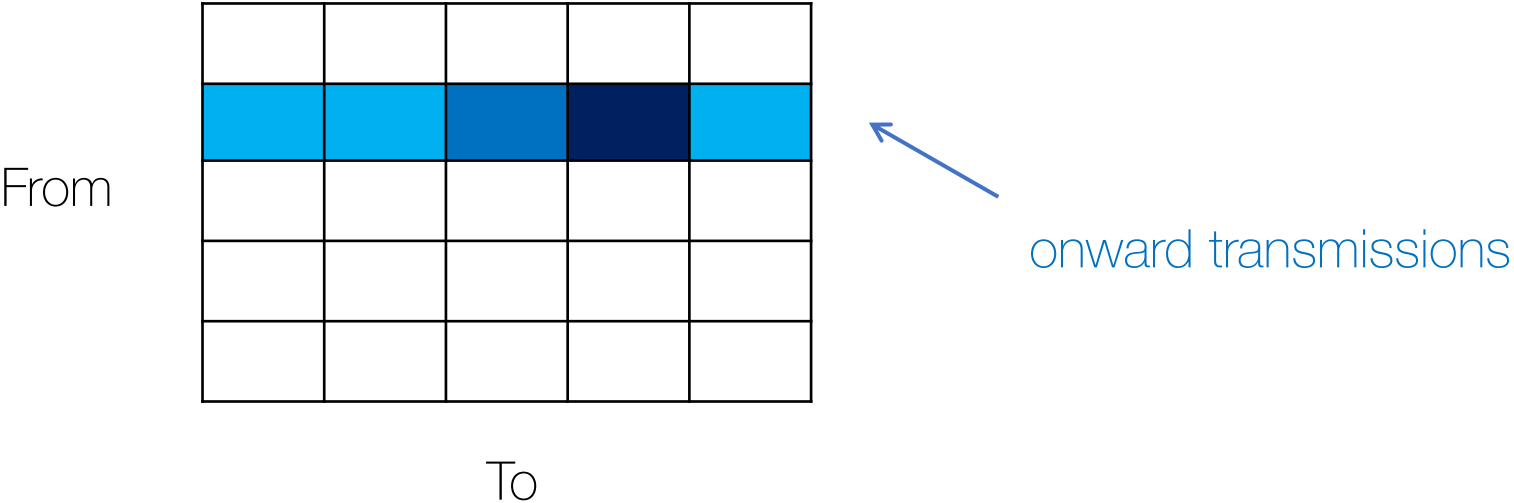
# Target quantities

- Target quantities derived from the transmission flow matrix:



# Target quantities

- Target quantities derived from the transmission flow matrix:



# Regression-type source attribution analysis



ORIGINAL ARTICLE | [Open Access](#) |

## Inferring the sources of HIV infection in Africa from deep-sequence data with semi-parametric Bayesian Poisson flow models

Xiaoyue Xi, Simon E. F. Spencer, Matthew Hall, M. Kate Grabowski, Joseph Kagaayi, Oliver Ratmann on behalf of [Rakai Health Sciences Program](#) and [PANGEA-HIV](#)

First published: 13 March 2022 | <https://doi.org/10.1111/rssc.12544>

Xi et al. JRSSC. 2022

Bu et al. – almost submitted.

Monod et al. – almost submitted.



# Regression-type source attribution analysis

$$Y_{p,i,j}^{g \rightarrow h} \sim \text{Poisson} \left( \xi_{p,j}^h \sum_{r \in p} \lambda_{r,i,j}^{g \rightarrow h} \right) \quad (6a)$$

$$\lambda_{r,i,j}^{g \rightarrow h} = \beta_{r,i,j}^{g \rightarrow h} \times S_{r,j}^h \times I_{r,i}^g \times |(t_r^{\text{end}} - t_r^{\text{start}})| \quad (6b)$$

$$\log \beta_{r,i,j}^{g \rightarrow h} = \hat{\mathbf{c}}^{g \rightarrow h}(i, j) + \gamma_0 + \gamma_g + \gamma_r + \gamma_{p(r)} + \mathbf{f}_0^{g \rightarrow h}(i, j) + \mathbf{f}_r^{g \rightarrow h}(j) + \mathbf{f}_{p(r)}^{g \rightarrow h}(i), \quad (6c)$$

$$\frac{\sum_i \lambda_{r,i,j}^{g \rightarrow h}}{S_{r,j}^h \times |(t_r^{\text{end}} - t_r^{\text{start}})|} \sim \text{LogNormal}$$

# Pros and cons

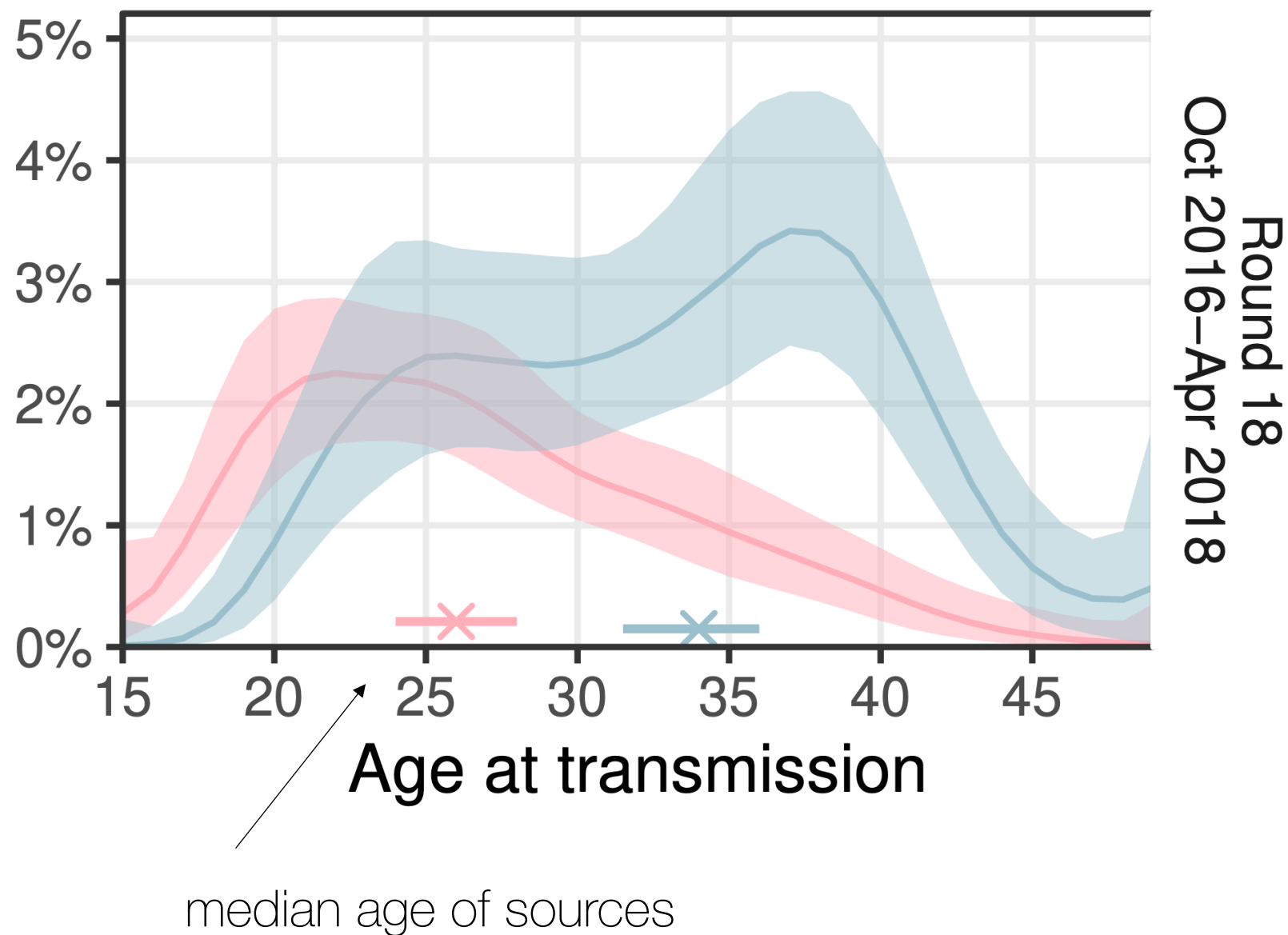
## Advantages:

- estimation of transmission flows computationally tractable (runtime several hours to 2-3 days)
- can be implemented in Stan
- can adjust flow estimates for observed sampling heterogeneity
- Gaussian process smoothing can be used to regularize inferences in highly-stratified populations

## Disadvantages:

- requires deep-sequence data
- does not use all available sequence data, only phylogenetically strongly supported source-recipient pairs

- **Age profile of male sources (blue), and female sources (pink)**
- **Blue + red = 100%**

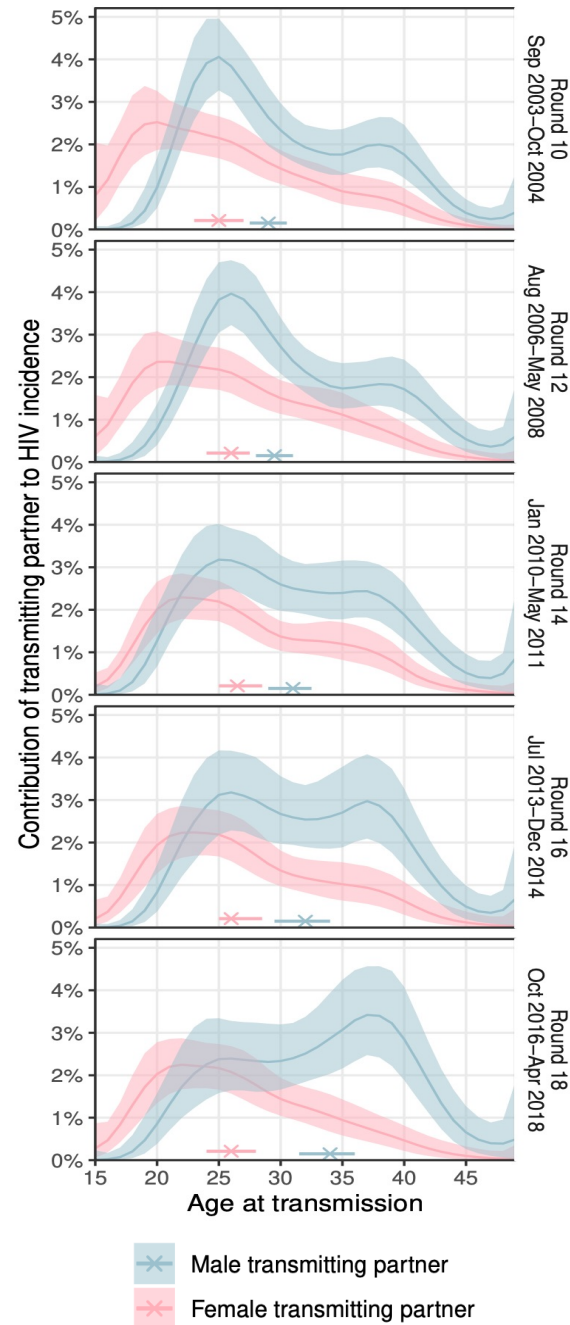


## %transmission from men

57.9%  
[56.1-59.6]

61.9%  
[60.2-63.7]

62.8%  
[60.2-65.2]



- **Proportion of transmissions from men is increasing**
- **Transmissions from men are shifting to older ages**
- **Disparities are widening, not closing**

# Changes in age/gender transmission flows

Transmission direction	Male-female difference in age at transmission	Infected partner by age at transmission			
		15-24 years (%) <sup>†</sup>	25-34 years (%) <sup>†</sup>	35-49 years (%) <sup>†</sup>	Total (%) <sup>†</sup>
Round 10, September 26, 2003 - November 23, 2004; 28 communities surveyed					
Male to female	Total	31.9% [30.2-33.6]	18.6% [17.7-19.6]	7.3% [6.7-7.9]	57.9% [56.1-59.6]
	<0 years	0.4% [0.2-0.7]	4.5% [3.0-6.3]	4.1% [2.7-5.5]	9.0% [6.8-11.5]
	0-6 years	16.0% [12.8-19.3]	8.5% [7.0-10.2]	3.0% [1.8-4.2]	27.5% [23.5-31.5]
	>6 years	15.5% [12.2-18.8]	5.6% [4.1-7.3]	0.2% [0.0-0.5]	21.3% [17.2-25.2]
Female to male	Total	14.8% [13.9-15.8]	20.7% [19.7-21.7]	6.6% [6.2-7.1]	42.1% [40.4-43.9]
	<0 years	6.8% [5.2-8.7]	4.3% [2.9-6.0]	0.4% [0.2-0.8]	11.6% [8.8-14.8]
	0-6 years	7.9% [5.9-9.9]	12.3% [10.5-13.9]	2.5% [1.7-3.3]	22.7% [19.7-25.7]
	>6 years	0.1% [0.0-0.2]	4.0% [2.7-5.8]	3.7% [2.7-4.7]	7.8% [5.8-10.1]
Total		46.7% [45.3-48.2]	39.3% [38.2-40.5]	13.9% [13.2-14.7]	100%
Round 15, August 10, 2011 - July 05, 2013; 33 communities surveyed					
Male to female	Total	32.2% [30.1-34.3]	22.0% [20.7-23.4]	7.7% [7.0-8.5]	61.9% [60.2-63.7]
	<0 years	0.5% [0.3-1.0]	4.8% [3.2-6.9]	3.9% [2.4-5.5]	9.3% [6.8-12.2]
	0-6 years	16.0% [12.7-19.4]	10.0% [8.1-12.0]	3.5% [2.1-4.9]	29.6% [25.3-33.9]
	>6 years	15.6% [12.2-19.1]	7.1% [5.3-9.1]	0.2% [0.1-0.7]	23.1% [18.6-27.3]
Female to male	Total	11.5% [10.6-12.4]	18.8% [17.8-19.9]	7.7% [7.1-8.4]	38.1% [36.3-39.8]
	<0 years	6.4% [4.9-7.8]	4.2% [2.9-5.9]	0.6% [0.2-1.2]	11.2% [8.7-14.0]
	0-6 years	5.1% [3.8-6.5]	11.8% [10.1-13.2]	3.2% [2.2-4.2]	20.0% [17.3-22.7]
	>6 years	0.0% [0.0-0.0]	2.8% [1.9-3.9]	3.9% [2.8-5.1]	6.8% [5.1-8.6]
Total		43.7% [41.9-45.6]	40.8% [39.3-42.4]	15.4% [14.6-16.4]	100%
Round 18, October 03, 2016 - May 22, 2018; 35 communities surveyed					
Male to female	Total	20.6% [18.2-23.4]	27.3% [25.3-29.4]	14.7% [13.3-16.3]	62.8% [60.2-65.2]
	<0 years	0.3% [0.1-0.7]	5.3% [3.2-8.4]	7.2% [4.8-9.7]	12.9% [9.2-17.4]
	0-6 years	8.7% [6.2-11.7]	13.2% [10.5-16.0]	7.0% [4.8-9.3]	29.0% [25.0-33.2]
	>6 years	11.5% [8.6-14.7]	8.6% [6.0-11.7]	0.5% [0.1-1.4]	20.7% [16.1-25.5]
Female to male	Total	11.2% [9.9-12.6]	17.3% [15.8-18.9]	8.7% [7.7-9.9]	37.2% [34.8-39.8]
	<0 years	5.8% [4.2-7.7]	3.5% [2.4-5.2]	0.4% [0.2-1.1]	9.8% [7.3-13.0]
	0-6 years	5.4% [3.6-7.1]	11.0% [9.3-12.6]	3.2% [2.1-4.5]	19.6% [16.7-22.4]
	>6 years	0.0% [0.0-0.1]	2.8% [1.8-4.0]	5.0% [3.5-6.4]	7.8% [5.8-9.9]
Total		31.9% [29.4-34.5]	44.7% [42.5-46.8]	23.4% [21.7-25.3]	100%

<sup>†</sup> Posterior median flow estimates and 95% credible intervals in each survey round.

# Transmission versus contacts

Monod et al. – almost submitted.

Chen et al. – in preparation.

Dan et al. arxiv 2022 <https://arxiv.org/pdf/2210.11358.pdf>



# Basic notations and definitions

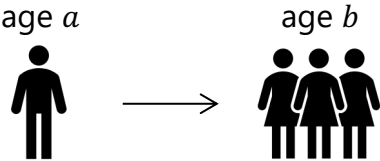
Participant's, contact's age / gender $a, b \in \{0, 1, \dots, 83, 84\}$ $g, h \in \{M, F\}$	Observed contact counts $Y_{ab}^{gh}$ $(Y_{ab}^{MF}, Y_{ab}^{FM}, Y_{ab}^{MM}, Y_{ab}^{FF})$	Age-gender-specific sample size $N_a^g$ $(N_a^M, N_a^F)$	Age-gender-specific pop. size $P_b^h$ $(P_b^M, P_b^F)$
---	--	---	---

# Basic notations and definitions

Participant's, contact's age / gender $a, b \in \{0, 1, \dots, 83, 84\}$ $g, h \in \{M, F\}$	Observed contact counts $Y_{ab}^{gh}$ $(Y_{ab}^{MF}, Y_{ab}^{FM}, Y_{ab}^{MM}, Y_{ab}^{FF})$	Age-gender-specific sample size $N_a^g$ $(N_a^M, N_a^F)$	Age-gender-specific pop. size $P_b^h$ $(P_b^M, P_b^F)$
---	--	---	---

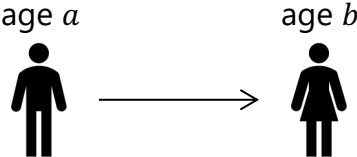
Contact intensity

$$m_{ab}^{MF} = \frac{\mathbb{E}[Y_{ab}^{MF}]}{N_a^M} = \frac{\mu_{ab}^{MF}}{N_a^M}$$



Contact rate

$$\gamma_{ab}^{MF} = \frac{m_{ab}^{MF}}{P_b^F} = \frac{\mu_{ab}^{MF}}{N_a^M P_b^F}$$

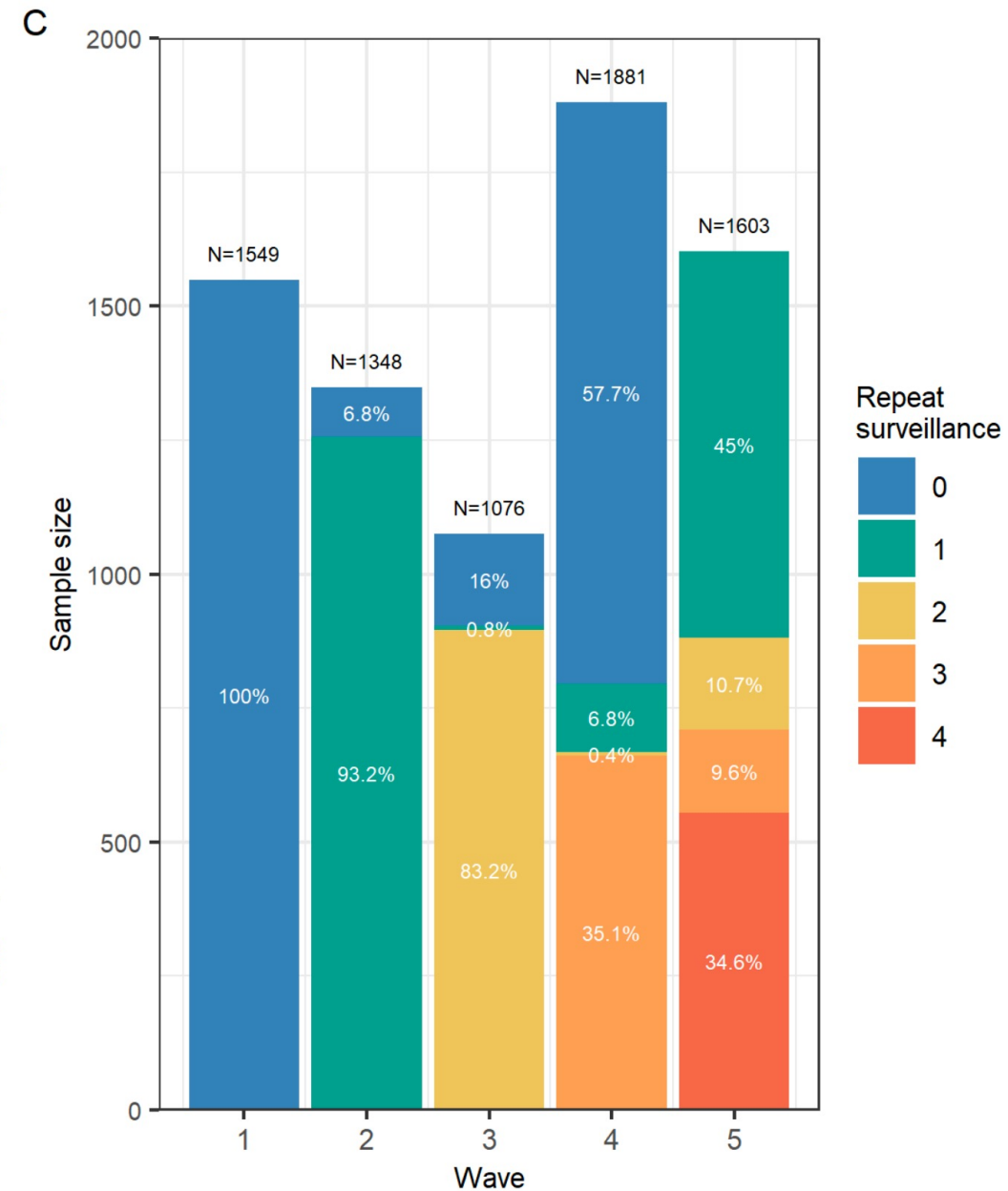
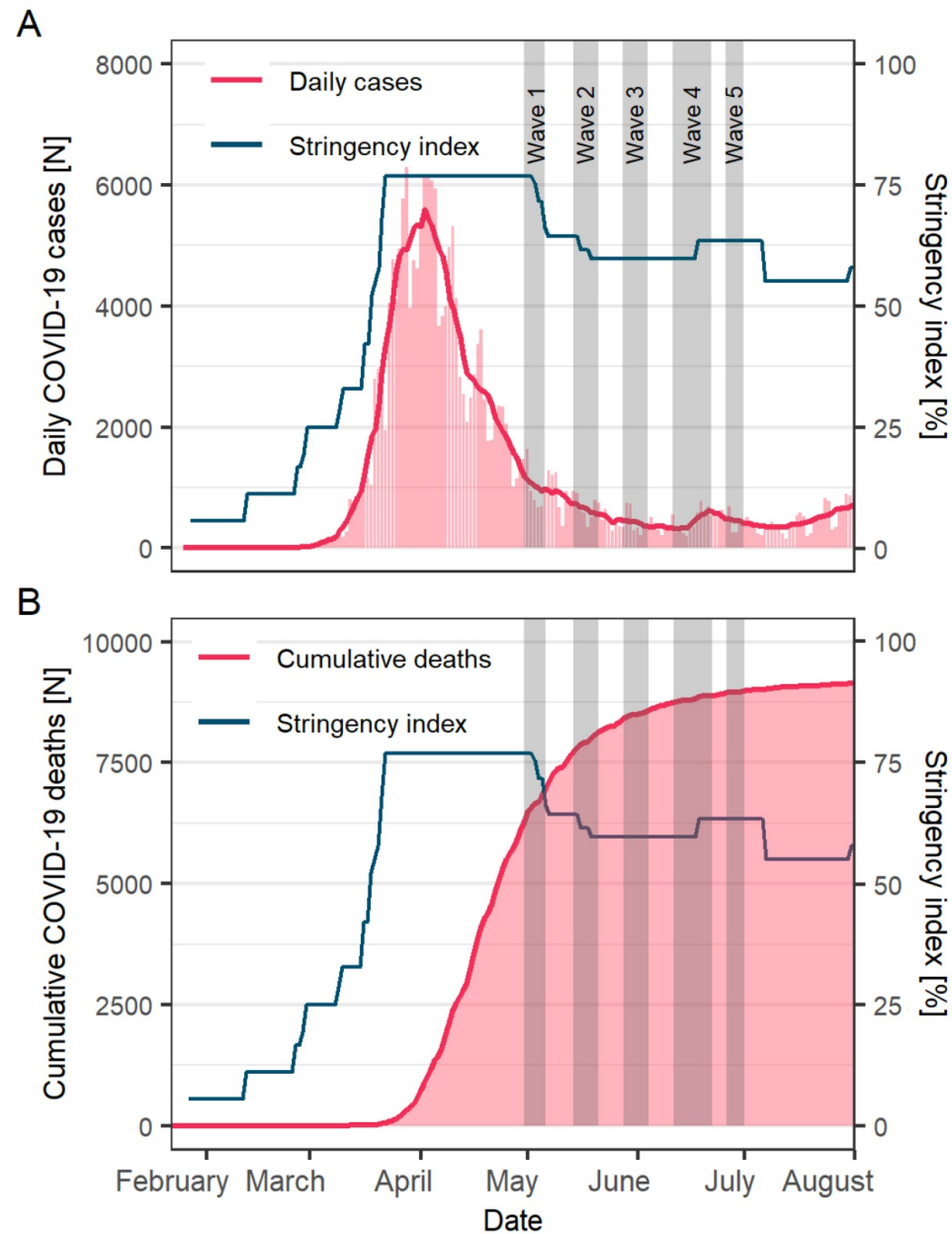


$$\gamma_{ab}^{MF} N_b^F P_a^M = \gamma_{ba}^{FM} N_a^F P_b^M$$

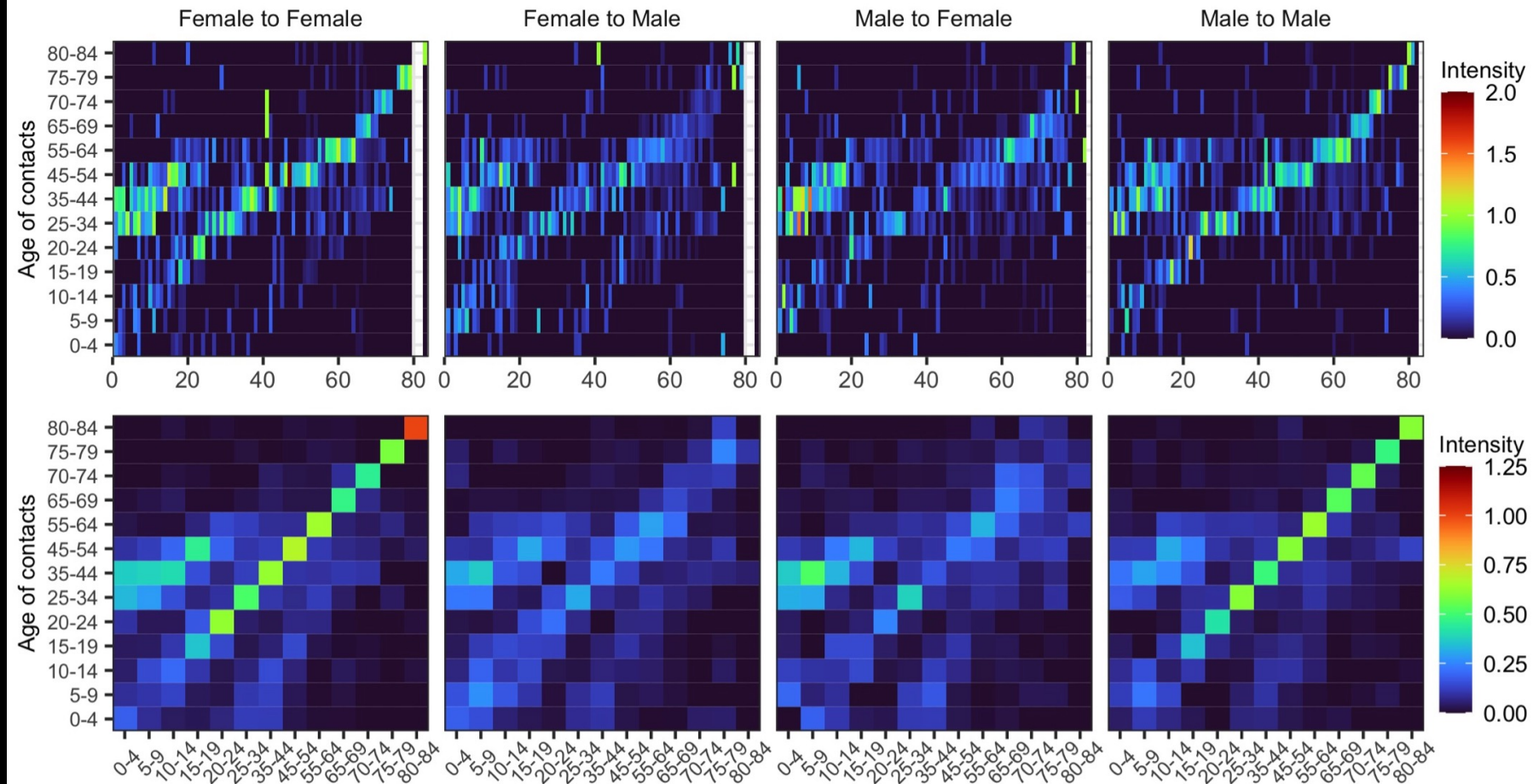
$$\gamma_{ab}^{MM} = \gamma_{ba}^{MM} \text{ for } a < b$$

$$\gamma_{ab}^{MF} = \gamma_{ba}^{FM}$$

# COVIMOD study, first 5 waves



# Data for wave 1, Germany



# Probabilistic modelling of contact patterns

**Model age-gender-specific contact counts with negative binomial**

$$Y_{ab}^{gh} \sim \text{NegBinomial}\left(\alpha_{ab}^{gh}, \frac{1-\nu}{\nu}\right)$$

$$\begin{aligned}\mu_{ab}^{gh} &= \alpha_{ab}^{gh} \frac{\nu}{1-\nu} \\ \log \mu_{ab}^{gh} &= \log m_{ab}^{gh} + \log N_a^g\end{aligned}$$

**Gender specific 2D offsets associated with 2D GP priors**

$$\log m_{ab}^{gh} = \beta_0 + f^{gh}(a, b) + \log P_b^h$$

**Exploit symmetry**

$$\begin{aligned}f^{MF}(a, b) &= f^{FM}(b, a) \quad \forall a, b \\ f^{MM}(a, b) &= f^{MM}(b, a) \quad a \leq b \\ f^{FF}(a, b) &= f^{FF}(b, a) \quad a \leq b\end{aligned}$$

# Recovering fine age structure from coarse age data

**Contacts' age is reported in discrete coarse age categories**

**Link our fine-age model to the coarse-age data by summing the shape parameter**

$$c \in \{0 - 4, 5 - 9, 10 - 14, \dots, 20 - 24, 25 - 34, 35 - 44, \dots, 65 - 69, 70 - 74, 80 - 84, \}$$

$$Y_{ac}^{gh} = \sum_{b \in c} Y_{ab}^{gh} \sim \text{NegBinomial} \left( \sum_{b \in c} \alpha_{ab}^{gh}, \frac{1 - \nu}{\nu} \right)$$



# Non-parametric modelling of contact rates

**Input is a 2D grid  $\mathbf{x}_1 = (\mathbf{a}_1, \mathbf{b}_1), \dots, \mathbf{x}_{AB} = (\mathbf{a}_A, \mathbf{b}_B)$ . A zero-mean multivariate Gaussian prior will have covariance matrix  $\mathbf{K} \in \mathbb{R}^{AB \times AB}$  with elements  $\mathbf{k}(\mathbf{x}_i, \mathbf{x}_j)$ .**

**Decompose.**

**Kronecker product.**

**Linear transformation of standard i.i.d. Gaussians.**

$$k((a, b), (a', b')) = k^1(a, a')k^2(b, b')$$

$$k^1(a, a') = \alpha^2 \exp\left(-\frac{(a - a')^2}{2\ell_a^2}\right)$$

$$k^2(b, b') = \alpha^2 \exp\left(-\frac{(b - b')^2}{2\ell_b^2}\right)$$

$$\mathbf{K} = \mathbf{K}^2 \otimes \mathbf{K}^1 = (\mathbf{L}^2 \otimes \mathbf{L}^1)(\mathbf{L}^2 \otimes \mathbf{L}^1)^T$$

$$f(x) = (\mathbf{L}^2 \otimes \mathbf{L}^1)z = \text{vec}\left(\left(\mathbf{L}^2(\mathbf{L}^1 \text{reshape}(z, A, B))^T\right)^T\right)$$

# Gaussian process approximations to alleviate computational bottleneck

**Cost of evaluating log-posterior for GP is  $\mathcal{O}(n^3)$ . Approximate the covariance kernel to reduce cost to  $\mathcal{O}(mn + m)$  where  $m \ll n$ .**

**Spectral density function of the covariance kernel + theory of pseudo-differential operators on compact space**

$$k^1(a, a') \approx \sum_{j=1}^{M^1} S^1 \left( \sqrt{\lambda_j^1} \right) \phi_j^1(a) \phi_j^1(a')$$

$$S^1(\omega) = \alpha^2 (2\pi \ell_a) \exp \left( -\frac{\ell_a^2 \omega}{2} \right)$$

$$\sqrt{\lambda_j^1} = \frac{j\pi}{2L^1}$$
$$\phi_j^1(x) = \sqrt{1/L^1} \sin \left( \sqrt{\lambda_j^1} (x + L^1) \right)$$

$$L^1 \approx \tilde{L}^1 = \Phi^1 \sqrt{\Delta^1}$$
$$f(x) = (\tilde{L}^2 \otimes \tilde{L}^1) \tilde{z}$$

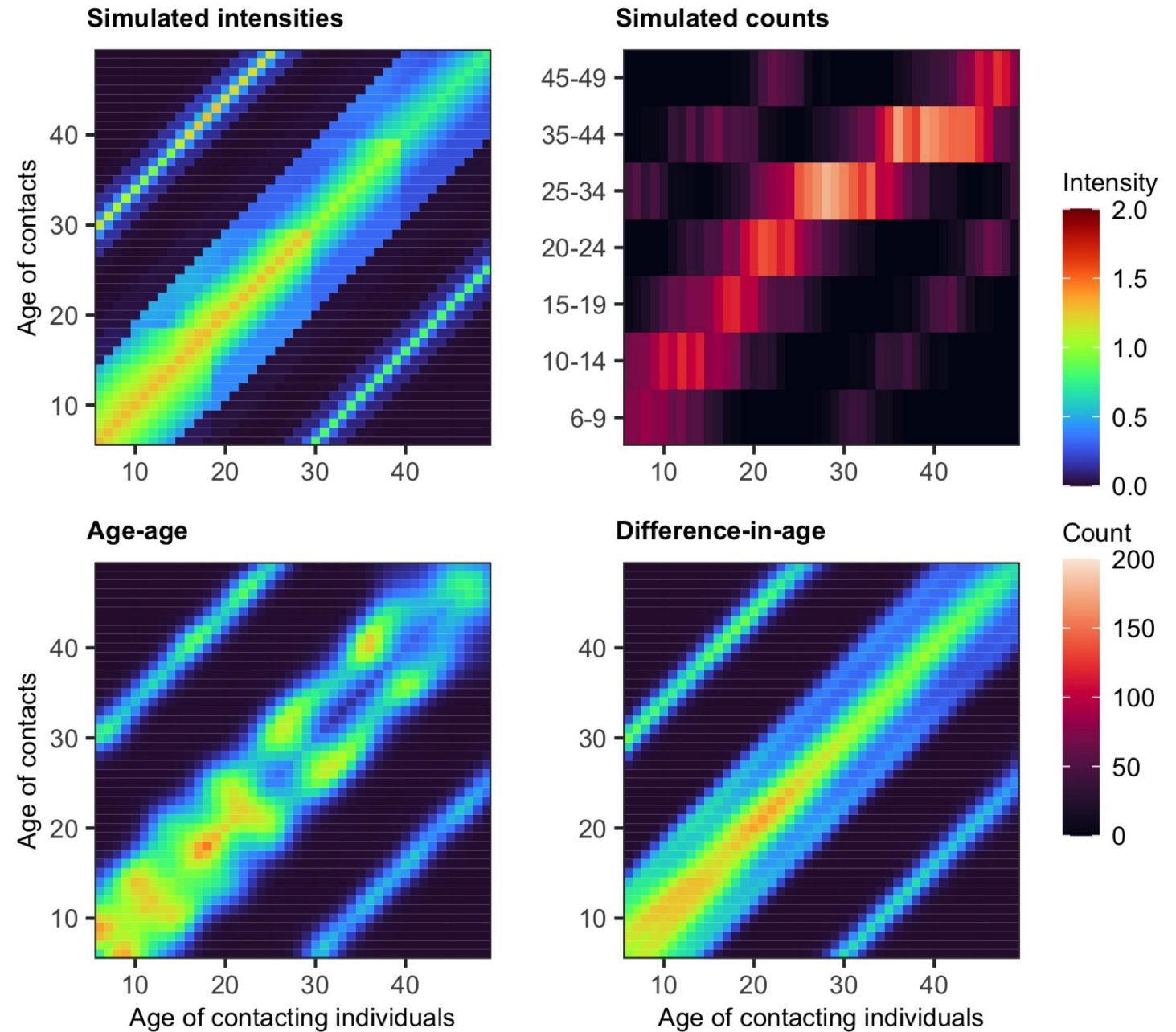
# Difference in age parameterization

**Human contact concentrate among individuals of similar age and individuals with similar age gaps**

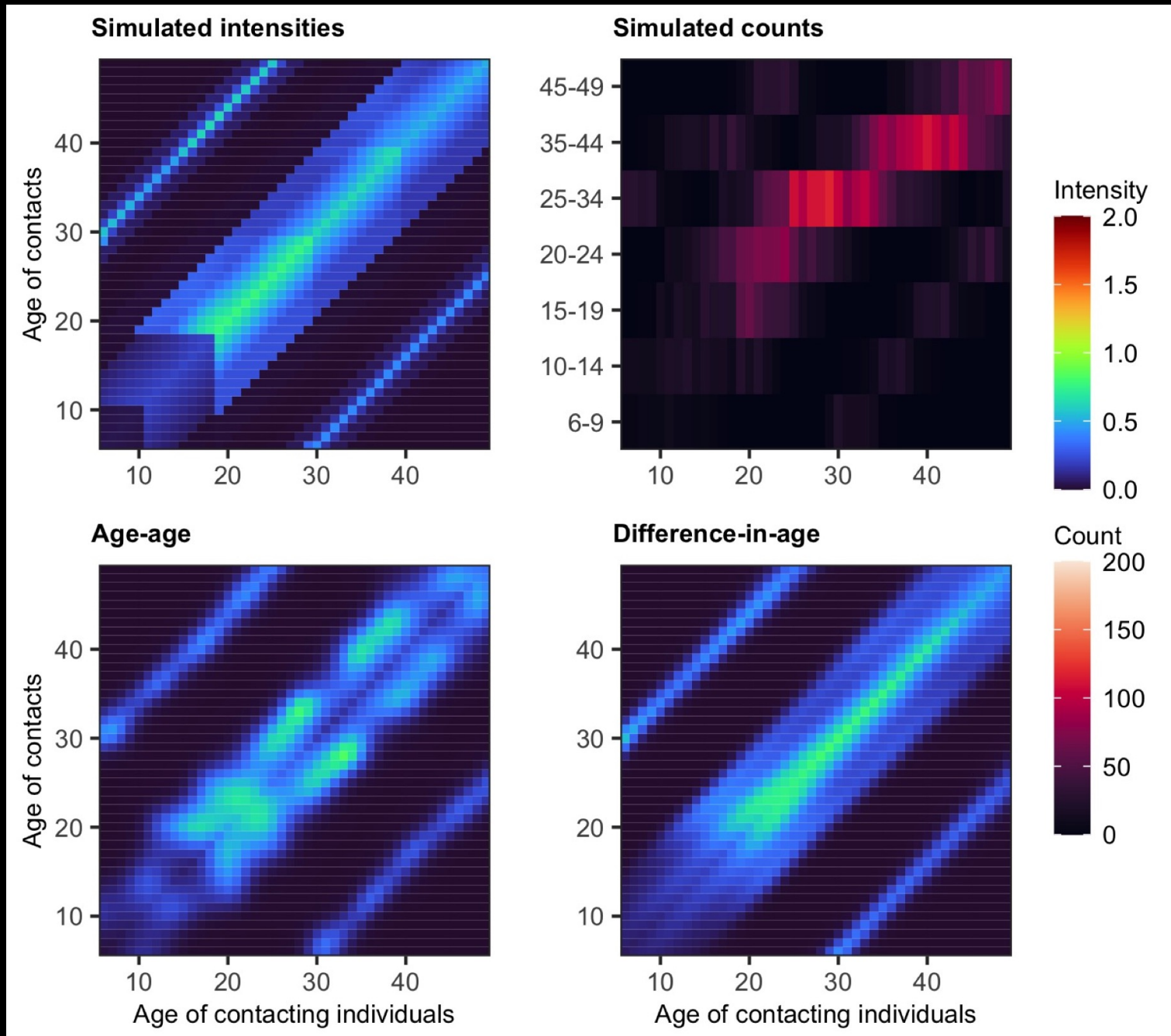
**Parameterize contact rate surface on a difference-in-age by age space as opposed to an age-by-age space**

$$\mathbf{f} = \begin{bmatrix} f_{11} & f_{12} & f_{13} & \dots & f_{1A} \\ f_{21} & f_{22} & f_{23} & & \\ f_{31} & f_{32} & f_{33} & & \\ \vdots & & & \ddots & \\ f_{A1} & & & & f_{AA} \end{bmatrix} \qquad \hat{\mathbf{f}} = \begin{bmatrix} & & & & f_{A1} \\ & & & \ddots & \vdots \\ & & f_{31} & \dots & f_{A,A-2} \\ & f_{21} & f_{32} & \dots & f_{A,A-1} \\ f_{11} & f_{22} & f_{33} & \dots & f_{AA} \\ f_{12} & f_{23} & \dots & \ddots & \\ f_{13} & \dots & \ddots & & \\ \vdots & \ddots & & & \\ f_{1A} & & & & \end{bmatrix}$$

# Accuracy on simulated data mimicking pre-COVID19 contact patterns

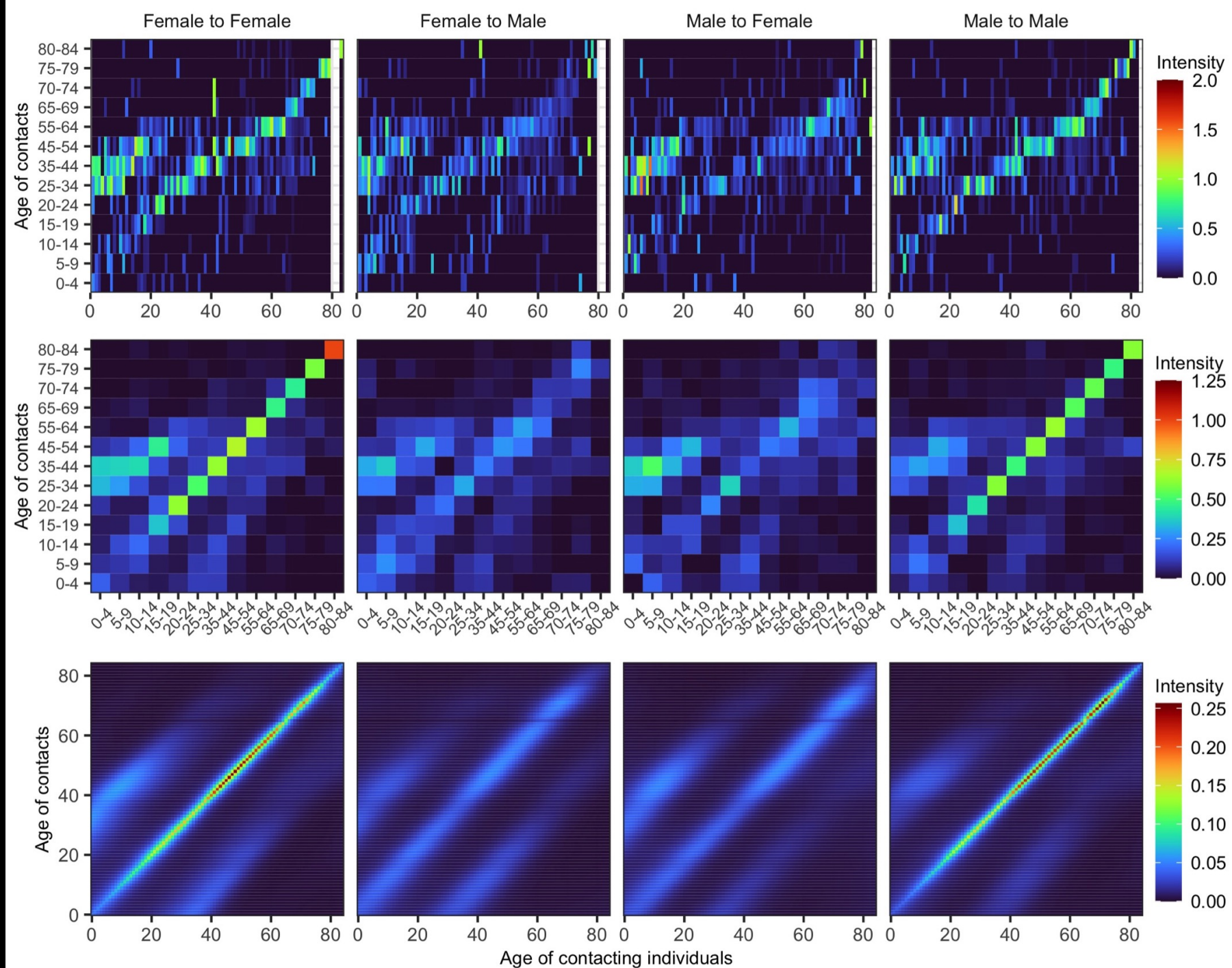


# Accuracy on simulated data mimicking in-COVID19 contact patterns



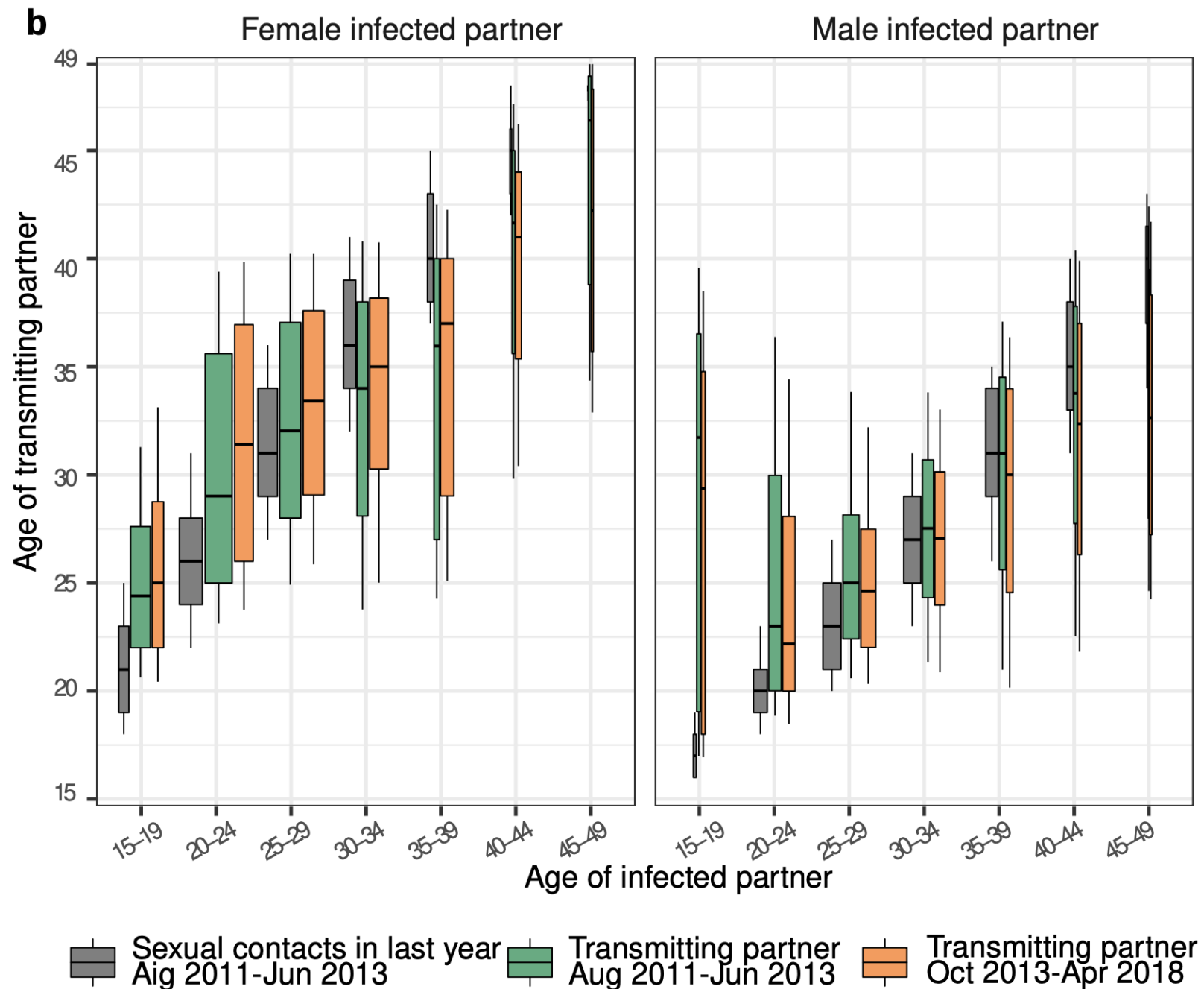


# Estimates for wave 1, Germany

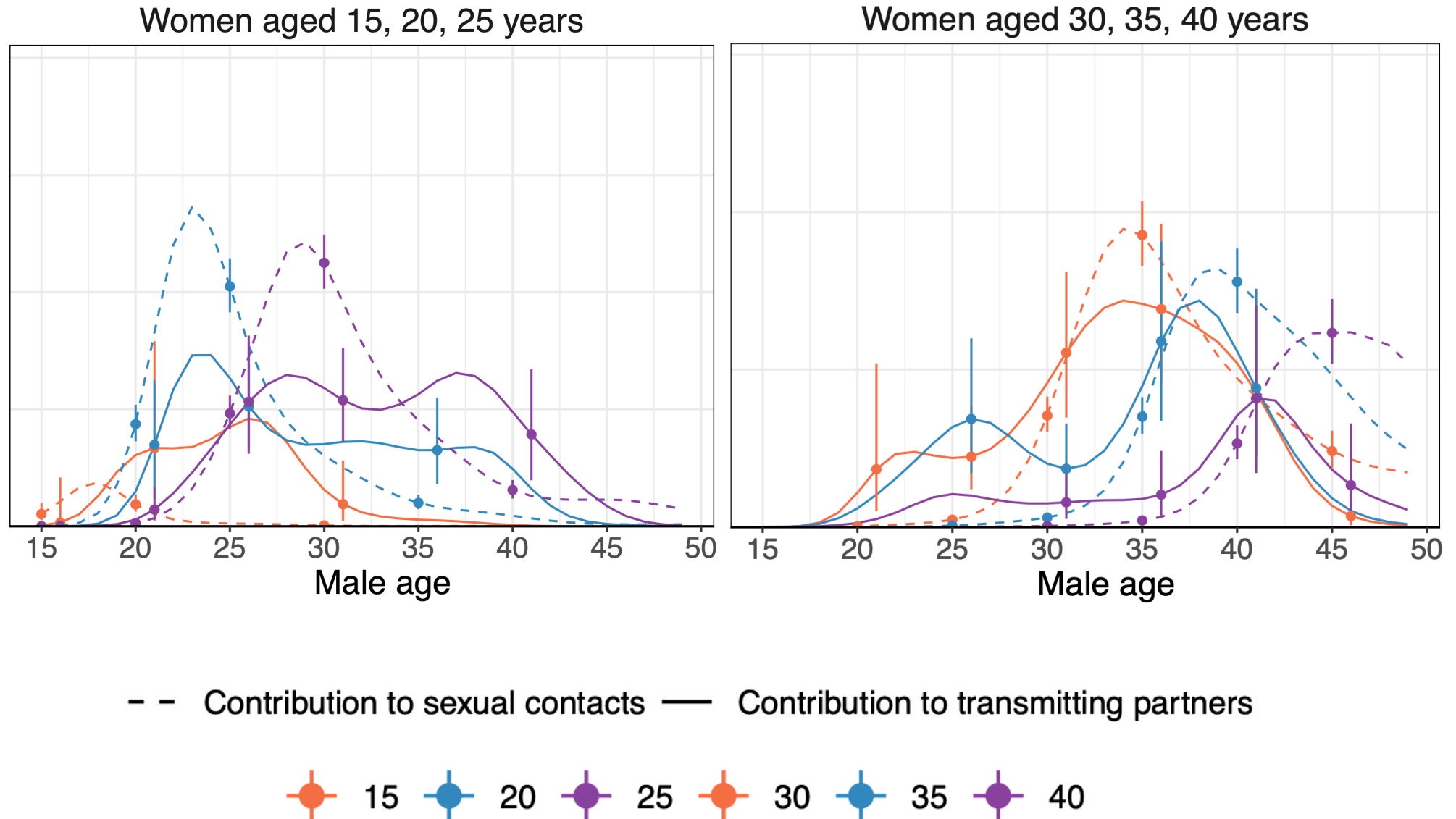




- **Adolescent girls and young women are infected by unusually older male partners.**
- **As women age, age difference between woman and infecting partner decreases.**



# Transmission flows vs. sexual contact patterns





- HIV incidence has declined faster among men than women.
- Average age of infection is increasing among women; and avg. age of transmission is increasing among men.
- While viral load suppression has increased in both genders, the viral load suppression gap has increased between men and women.
- Men are accounting for an increasing proportion of transmissions.
- Having closed the viral load suppression gap between men in women, would have reduced female HIV incidence by 50%.





Thank you

# Acknowledgments

## Rakai Health Sciences Program

David Serwadda  
Fred Nalugoda  
Joseph Kagaayi  
Godfrey Kigozi  
Gertrude Nakigozi  
Tom Lutalo  
Robert Ssekubugu  
Grace Kigozi  
Jeremiah Bazaale  
Edward Kankaka  
Ronald Galiwango  
Victor Ssempijja

## Johns Hopkins Bloomberg School of Public Health

Mary Kate Grabowski  
Ronald Gray  
Maria Wawer  
Caitlin Kennedy  
Joseph Ssekasanvu

## Imperial College London

Melodie Monod  
Alexandra Blenkinsop  
Andrea Brizzi  
Yu Chen  
Xiaoyue Xi  
Shozen Dan

## Johns Hopkins School of Medicine

Aaron Tobian  
Larry W Chang

## National Institute of Allergy and Infectious Diseases

Thomas Quinn  
Andrew Redd  
Steven J Reynolds  
Oliver Laeyendecker

## Oxford University

Christophe Fraser  
Matthew Hall  
Chris Wymant  
Tanya Golubchik  
Lucie Abeler-Dorner  
David Bonsall  
Laura Thompson

## LSHTM

Peter Godfrey-Fausset

## Institute for Disease Modelling

Adam Akullian

## University of Warwick

Simon Spencer

## Rakai Health Science Program Staff and Study participants



BILL & MELINDA  
GATES *foundation*

