Sheffield Hallam University

Mapping age- and sex-specific HIV prevalence in adults in sub-Saharan Africa, 2000–2018

HAEUSER, Emily, SERFES, Audrey L., CORK, Michael A., YANG, Mingyou, ABBASTABAR, Hedayat, ABHILASH, E. S., ADABI, Maryam, ADEBAYO, Oladimeji M., ADEKANMBI, Victor, ADEYINKA, Daniel Adedavo, AFZAL, Saira, AHINKORAH, Bright Opoku, AHMADI, Keivan, AHMED, Muktar Beshir, AKALU, Yonas, AKINYEMI, Rufus Olusola, AKUNNA, Chisom Joyqueenet, ALAHDAB, Fares, ALANEZI, Fahad Mashhour, ALANZI, Turki M., ALENE, Kefyalew Addis, ALHASSAN, Robert Kaba, ALIPOUR, Vahid, ALMASI-HASHIANI, Amir, ALVIS-GUZMAN, Nelson, AMEYAW, Edward Kwabena, AMINI, Saeed, AMUGSI, Dickson A., ANCUCEANU, Robert, ANVARI, Davood, APPIAH, Seth Christopher Yaw, ARABLOO, Jalal, AREMU, Olatunde, ASEMAHAGN, Mulusew A., JAFARABADI, Mohammad Asghari, AWEDEW, Atalel Fentahun, QUINTANILLA, Beatriz Paulina Ayala, AYANORE, Martin Amogre, AYNALEM, Yared Asmare, AZARI, Samad, AZENE, Zelalem Nigussie, DARSHAN, B. B., BABALOLA, Tesleem Kayode, BAIG, Atif Amin, BANACH, Maciej, BÄRNIGHAUSEN, Till Winfried, BELL, Arielle Wilder, BHAGAVATHULA, Akshaya Srikanth, BHARDWAJ, Nikha, BHARDWAJ, Pankaj, BHATTACHARYYA, Krittika, BIJANI, Ali, BITEW, Zebenay Workneh, BOHLOULI, Somayeh, BOLARINWA, Obasanjo Afolabi, BOLOOR, Archith, BOZICEVIC, Ivana, BUTT, Zahid A., CÁRDENAS, Rosario, CARVALHO, Felix, CHARAN, Jaykaran, CHATTU, Vijay Kumar, CHOWDHURY, Mohiuddin Ahsanul Kabir, CHU, Dinh-Toi, COWDEN, Richard G., DAHLAWI, Saad M. A., DAMIANI, Giovanni, DARTEH, Eugene Kofuor Maafo, DARWESH, Aso Mohammad, DAS NEVES, José, WEAVER, Nicole Davis, DE LEO, Diego, DE NEVE, Jan-Walter, DERIBE, Kebede, DEUBA, Keshab, DHARMARATNE, Samath, DIANATINASAB, Mostafa, DIAZ, Daniel, DIDARLOO, Alireza, DJALALINIA, Shirin, DOROSTKAR, Fariba, DUBLJANIN, Eleonora, DUKO, Bereket, EL TANTAWI, Maha, EL-JAAFARY, Shaimaa I., ESHRATI, Babak, ESKANDARIEH, Sharareh, EYAWO, Oghenowede, EZEONWUMELU, Ifeanyi Jude, EZZIKOURI, Sayeh, FARZADFAR, Farshad, FATTAHI, Nazir, FAUK, Nelsensius Klau, FERNANDES, Eduarda, FILIP, Irina, FISCHER, Florian, FOIGT, Nataliya A., FOROUTAN, Masoud, FUKUMOTO, Takeshi, GAD, Mohamed M., GAIDHANE, Abhay Motiramji, GEBREGIORGIS, Birhan Gebresillassie, GEBREMEDHIN, Ketema Bizuwork, GETACHER, Lemma, GHADIRI, Keyghobad, GHASHGHAEE, Ahmad, GOLECHHA, Mahaveer, GUBARI, Mohammed Ibrahim Mohialdeen, GUGNANI, Harish Chander, GUIMARÃES, Rafael Alves, HAIDER, Mohammad Rifat, HAJ-MIRZAIAN, Arvin, HAMIDI, Samer, HASHI, Abdiwahab, HASSANIPOUR, Soheil, HASSANKHANI, Hadi, HAYAT, Khezar, HERTELIU, Claudiu, HO, Hung Chak,

HOLLA, Ramesh, HOSSEINI, Mostafa, HOSSEINZADEH, Mehdi, HWANG, Bing-Fang, IBITOYE, Segun Emmanuel, ILESANMI, Olayinka Stephen, ILIC, Irena M., ILIC, Milena D., ISLAM, Rakibul M., IWU, Chidozie C. D., JAKOVLJEVIC, Mihajlo, JHA, Ravi Prakash, JI, John S., JOHNSON, Kimberly B., JOSEPH, Nitin, JOSHUA, Vasna, JOUKAR, Farahnaz, JOZWIAK, Jacek Jerzy, KALANKESH, Leila R., KALHOR, Rohollah, KAMYARI, Naser, KANCHAN, Tanuj, MATIN, Behzad Karami, KARIMI, Salah Eddin, KAYODE, Gbenga A., KARYANI, Ali Kazemi, KERAMATI, Maryam, KHAN, Ejaz Ahmad, KHAN, Gulfaraz, KHAN, Md Nuruzzaman, KHATAB, Khaled http://orcid.org/0000-0002-8755-3964>, KHUBCHANDANI, Jagdish, KIM, Yun Jin, KISA, Adnan, KISA, Sezer, KOPEC, Jacek A., KOSEN, Soewarta, LAXMINARAYANA, Sindhura Lakshmi Koulmane, KOYANAGI, Ai, KRISHAN, Kewal, DEFO, Barthelemy Kuate, KUGBEY, Nuworza, KULKARNI, Vaman, KUMAR, Manasi, KUMAR, Nithin, KUSUMA, Dian, LA VECCHIA, Carlo, LAL, Dharmesh Kumar, LANDIRES, Iván, LARSON, Heidi Jane, LASRADO, Savita, LEE, Paul H., LI, Shanshan, LIU, Xuefeng, MALEKI, Afshin, MALIK, Preeti, MANSOURNIA, Mohammad Ali, MARTINS-MELO, Francisco Rogerlândio, MENDOZA, Walter, MENEZES, Ritesh G., MENGESHA, Endalkachew Worku, MERETOJA, Tuomo J., MESTROVIC, Tomislav, MIRICA, Andreea, MOAZEN, Babak, MOHAMAD, Osama, MOHAMMAD, Yousef, MOHAMMADIAN-HAFSHEJANI, Abdollah, MOHAMMADPOURHODKI, Reza, MOHAMMED, Salahuddin, MOHAMMED, Shafiu, MOKDAD, Ali H., MORADI, Masoud, MORAGA, Paula, MUBARIK, Sumaira, MULU, Getaneh Baye B., MWANRI, Lillian, NAGARAJAN, Ahamarshan Jayaraman, NAIMZADA, Mukhammad David, NAVEED, Muhammad, NAZARI, Javad, NDEJJO, Rawlance, NEGOI, Ionut, NGALESONI, Frida N., NGUEFACK-TSAGUE, Georges, NGUNJIRI, Josephine W., NGUYEN, Cuong Tat, NGUYEN, Huong Lan Thi, NNAJI, Chukwudi A., NOUBIAP, Jean Jacques, NUÑEZ-SAMUDIO, Virginia, NWATAH, Vincent Ebuka, OANCEA, Bogdan, ODUKOYA, Oluwakemi Ololade, OLAGUNJU, Andrew T., OLAKUNDE, Babayemi Oluwaseun, OLUSANYA, Bolajoko Olubukunola, OLUSANYA, Jacob Olusegun, BALI, Ahmed Omar, ONWUJEKWE, Obinna E., ORISAKWE, Orish Ebere, OTSTAVNOV, Nikita, OTSTAVNOV, Stanislav S., OWOLABI, Mayowa O., MAHESH, P. A., PADUBIDRI, Jagadish Rao, PANA, Adrian, PANDEY, Ashok, PANDI-PERUMAL, Seithikurippu R., KAN, Fatemeh Pashazadeh, PATTON, George C., PAWAR, Shrikant, PEPRAH, Emmanuel K., POSTMA, Maarten J., PREOTESCU, Liliana, SYED, Zahiruddin Quazi, RABIEE, Navid, RADFAR, Amir, RAFIEI, Alireza, RAHIM, Fakher, RAHIMI-MOVAGHAR, Vafa, RAHMANI, Amir Masoud, RAMEZANZADEH, Kiana, RANA, Juwel, RANABHAT, Chhabi Lal, RAO, Sowmya J., RAWAF, David Laith, RAWAF, Salman, RAWASSIZADEH, Reza, REGASSA, Lemma Demissie, REZAEI, Nima, REZAPOUR, Aziz, RIAZ, Mavra A., RIBEIRO, Ana Isabel, ROSS, Jennifer M., RUBAGOTTI, Enrico, RUMISHA, Susan Fred, RWEGERERA, Godfrey M., MOGHADDAM, Sahar Saeedi, SAGAR, Rajesh, SAHILEDENGLE, Biniyam, SAHU, Maitreyi, SALEM, Marwa Rashad, KAFIL, Hossein Samadi, SAMY, Abdallah M., SARTORIUS, Benn, SATHIAN, Brijesh, SEIDU, Abdul-Aziz, SHAHEEN, Amira A., SHAIKH, Masood Ali, SHAMSIZADEH, Morteza, SHIFERAW, Wondimeneh Shibabaw, SHIN, Jae II,

SHRESTHA, Roman, SINGH, Jasvinder A., SKRYABIN, Valentin Yurievich, SKRYABINA, Anna Aleksandrovna, SOLTANI, Shahin, SUFIYAN, Mu'awiyyah Babale, TABUCHI, Takahiro, TADESSE, Eyayou Girma, TAVEIRA, Nuno, TESFAY, Fisaha Haile, THAPAR, Rekha, TOVANI-PALONE, Marcos Roberto, TSEGAYE, Gebiyaw Wudie, UMEOKONKWO, Chukwuma David, UNNIKRISHNAN, Bhaskaran, VILLAFAÑE, Jorge Hugo, VIOLANTE, Francesco S., VO, Bay, VU, Giang Thu, WADO, Yohannes Dibaba, WAHEED, Yasir, WAMAI, Richard G., WANG, Yanzhong, WARD, Paul, WICKRAMASINGHE, Nuwan Darshana, WILSON, Katherine, YAYA, Sanni, YIP, Paul, YONEMOTO, Naohiro, YU, Chuanhua, ZASTROZHIN, Mikhail Sergeevich, ZHANG, Yunquan, ZHANG, Zhi-Jiang, HAY, Simon I. and DWYER-LINDGREN, Laura

Available from Sheffield Hallam University Research Archive (SHURA) at:

https://shura.shu.ac.uk/31185/

This document is the Supplemental Material

Citation:

HAEUSER, Emily, SERFES, Audrey L., CORK, Michael A., YANG, Mingyou, ABBASTABAR, Hedayat, ABHILASH, E. S., ADABI, Maryam, ADEBAYO, Oladimeji M., ADEKANMBI, Victor, ADEYINKA, Daniel Adedayo, AFZAL, Saira, AHINKORAH, Bright Opoku, AHMADI, Keivan, AHMED, Muktar Beshir, AKALU, Yonas, AKINYEMI, Rufus Olusola, AKUNNA, Chisom Joygueenet, ALAHDAB, Fares, ALANEZI, Fahad Mashhour, ALANZI, Turki M., ALENE, Kefyalew Addis, ALHASSAN, Robert Kaba, ALIPOUR, Vahid, ALMASI-HASHIANI, Amir, ALVIS-GUZMAN, Nelson, AMEYAW, Edward Kwabena, AMINI, Saeed, AMUGSI, Dickson A., ANCUCEANU, Robert, ANVARI, Davood, APPIAH, Seth Christopher Yaw, ARABLOO, Jalal, AREMU, Olatunde, ASEMAHAGN, Mulusew A., JAFARABADI, Mohammad Asghari, AWEDEW, Atalel Fentahun, OUINTANILLA, Beatriz Paulina Ayala, AYANORE, Martin Amogre, AYNALEM, Yared Asmare, AZARI, Samad, AZENE, Zelalem Nigussie, DARSHAN, B. B., BABALOLA, Tesleem Kayode, BAIG, Atif Amin, BANACH, Maciej, BÄRNIGHAUSEN, Till Winfried, BELL, Arielle Wilder, BHAGAVATHULA, Akshava Srikanth, BHARDWAJ, Nikha, BHARDWAJ, Pankai, BHATTACHARYYA, Krittika, BIJANI, Ali, BITEW, Zebenay Workneh, BOHLOULI, Somayeh, BOLARINWA, Obasanjo Afolabi, BOLOOR, Archith, BOZICEVIC, Ivana, BUTT, Zahid A., CÁRDENAS, Rosario, CARVALHO, Felix, CHARAN, Jaykaran, CHATTU, Vijay Kumar, CHOWDHURY, Mohiuddin Ahsanul Kabir, CHU, Dinh-Toi, COWDEN, Richard G., DAHLAWI, Saad M. A., DAMIANI, Giovanni, DARTEH, Eugene Kofuor Maafo, DARWESH, Aso Mohammad, DAS NEVES, José, WEAVER, Nicole Davis, DE LEO, Diego, DE NEVE, Jan-Walter, DERIBE, Kebede, DEUBA, Keshab, DHARMARATNE, Samath, DIANATINASAB, Mostafa, DIAZ, Daniel, DIDARLOO, Alireza, DJALALINIA, Shirin, DOROSTKAR, Fariba, DUBLJANIN,

Eleonora, DUKO, Bereket, EL TANTAWI, Maha, EL-JAAFARY, Shaimaa I., ESHRATI, Babak, ESKANDARIEH, Sharareh, EYAWO, Oghenowede, EZEONWUMELU, Ifeanyi Jude, EZZIKOURI, Sayeh, FARZADFAR, Farshad, FATTAHI, Nazir, FAUK, Nelsensius Klau, FERNANDES, Eduarda, FILIP, Irina, FISCHER, Florian, FOIGT, Nataliya A., FOROUTAN, Masoud, FUKUMOTO, Takeshi, GAD, Mohamed M., GAIDHANE, Abhay Motiramji, GEBREGIORGIS, Birhan Gebresillassie, GEBREMEDHIN, Ketema Bizuwork, GETACHER, Lemma, GHADIRI, Kevghobad, GHASHGHAEE, Ahmad, GOLECHHA, Mahaveer, GUBARI, Mohammed Ibrahim Mohialdeen, GUGNANI, Harish Chander, GUIMARÃES, Rafael Alves, HAIDER, Mohammad Rifat, HAJ-MIRZAIAN, Arvin, HAMIDI, Samer, HASHI, Abdiwahab, HASSANIPOUR, Soheil, HASSANKHANI, Hadi, HAYAT, Khezar, HERTELIU, Claudiu, HO, Hung Chak, HOLLA, Ramesh, HOSSEINI, Mostafa, HOSSEINZADEH, Mehdi, HWANG, Bing-Fang, IBITOYE, Segun Emmanuel, ILESANMI, Olayinka Stephen, ILIC, Irena M., ILIC, Milena D., ISLAM, Rakibul M., IWU. Chidozie C. D., JAKOVLJEVIC, Mihajlo, JHA, Ravi Prakash, JI, John S., JOHNSON, Kimberly B., JOSEPH, Nitin, JOSHUA, Vasna, JOUKAR, Farahnaz, JOZWIAK, Jacek Jerzy, KALANKESH, Leila R., KALHOR, Rohollah, KAMYARI, Naser, KANCHAN, Tanuj, MATIN, Behzad Karami, KARIMI, Salah Eddin, KAYODE, Gbenga A., KARYANI, Ali Kazemi, KERAMATI, Maryam, KHAN, Ejaz Ahmad, KHAN, Gulfaraz, KHAN, Md Nuruzzaman, KHATAB, Khaled, KHUBCHANDANI, Jagdish, KIM, Yun Jin, KISA, Adnan, KISA, Sezer, KOPEC, Jacek A., KOSEN, Soewarta, LAXMINARAYANA, Sindhura Lakshmi Koulmane, KOYANAGI, Ai, KRISHAN, Kewal, DEFO, Barthelemy Kuate, KUGBEY, Nuworza, KULKARNI, Vaman, KUMAR, Manasi, KUMAR, Nithin, KUSUMA, Dian, LA VECCHIA, Carlo, LAL, Dharmesh Kumar, LANDIRES, Iván, LARSON, Heidi Jane, LASRADO, Savita, LEE, Paul H., LI, Shanshan, LIU, Xuefeng, MALEKI, Afshin, MALIK, Preeti, MANSOURNIA, Mohammad Ali, MARTINS-MELO, Francisco Rogerlândio, MENDOZA, Walter, MENEZES, Ritesh G., MENGESHA, Endalkachew Worku, MERETOJA, Tuomo J., MESTROVIC, Tomislav, MIRICA, Andreea, MOAZEN, Babak, MOHAMAD, Osama, MOHAMMAD. Yousef. MOHAMMADIAN-HAFSHEJANI. Abdollah. MOHAMMADPOURHODKI, Reza, MOHAMMED, Salahuddin, MOHAMMED, Shafiu, MOKDAD, Ali H., MORADI, Masoud, MORAGA, Paula, MUBARIK, Sumaira, MULU, Getaneh Baye B., MWANRI, Lillian, NAGARAJAN, Ahamarshan Jayaraman, NAIMZADA, Mukhammad David, NAVEED, Muhammad, NAZARI, Javad, NDEJJO, Rawlance, NEGOI, Ionut, NGALESONI, Frida N., NGUEFACK-TSAGUE, Georges, NGUNJIRI, Josephine W., NGUYEN, Cuong Tat, NGUYEN, Huong Lan Thi, NNAJI, Chukwudi A., NOUBIAP, Jean Jacques, NUÑEZ-SAMUDIO, Virginia, NWATAH, Vincent Ebuka, OANCEA, Bogdan, ODUKOYA, Oluwakemi Ololade, OLAGUNJU, Andrew T., OLAKUNDE, Babayemi Oluwaseun, OLUSANYA, Bolajoko Olubukunola, OLUSANYA, Jacob Olusegun, BALI, Ahmed Omar, ONWUJEKWE, Obinna E., ORISAKWE, Orish Ebere, OTSTAVNOV, Nikita, OTSTAVNOV, Stanislav S., OWOLABI, Mayowa O., MAHESH, P. A., PADUBIDRI, Jagadish Rao, PANA, Adrian, PANDEY, Ashok, PANDI-PERUMAL, Seithikurippu R., KAN, Fatemeh Pashazadeh, PATTON, George C., PAWAR, Shrikant, PEPRAH, Emmanuel K., POSTMA, Maarten J., PREOTESCU, Liliana, SYED, Zahiruddin Quazi, RABIEE, Navid, RADFAR, Amir, RAFIEI, Alireza, RAHIM, Fakher, RAHIMI-MOVAGHAR, Vafa, RAHMANI, Amir Masoud, RAMEZANZADEH, Kiana, RANA, Juwel, RANABHAT, Chhabi Lal, RAO, Sowmya J., RAWAF, David Laith, RAWAF, Salman, RAWASSIZADEH, Reza, REGASSA, Lemma Demissie, REZAEI, Nima, REZAPOUR, Aziz, RIAZ, Mavra A., RIBEIRO, Ana Isabel, ROSS, Jennifer M., RUBAGOTTI, Enrico, RUMISHA, Susan Fred, RWEGERERA, Godfrey M., MOGHADDAM, Sahar Saeedi, SAGAR, Rajesh, SAHILEDENGLE, Biniyam, SAHU, Maitreyi, SALEM, Marwa Rashad, KAFIL,

Hossein Samadi, SAMY, Abdallah M., SARTORIUS, Benn, SATHIAN, Brijesh, SEIDU, Abdul-Aziz, SHAHEEN, Amira A., SHAIKH, Masood Ali, SHAMSIZADEH, Morteza, SHIFERAW, Wondimeneh Shibabaw, SHIN, Jae II, SHRESTHA, Roman, SINGH, Jasvinder A., SKRYABIN, Valentin Yurievich, SKRYABINA, Anna Aleksandrovna, SOLTANI, Shahin, SUFIYAN, Mu'awiyyah Babale, TABUCHI, Takahiro, TADESSE, Eyayou Girma, TAVEIRA, Nuno, TESFAY, Fisaha Haile, THAPAR, Rekha, TOVANI-PALONE, Marcos Roberto, TSEGAYE, Gebiyaw Wudie, UMEOKONKWO, Chukwuma David, UNNIKRISHNAN, Bhaskaran, VILLAFAÑE, Jorge Hugo, VIOLANTE, Francesco S., VO, Bay, VU, Giang Thu, WADO, Yohannes Dibaba, WAHEED, Yasir, WAMAI, Richard G., WANG, Yanzhong, WARD, Paul, WICKRAMASINGHE, Nuwan Darshana, WILSON, Katherine, YAYA, Sanni, YIP, Paul, YONEMOTO, Naohiro, YU, Chuanhua, ZASTROZHIN, Mikhail Sergeevich, ZHANG, Yunquan, ZHANG, Zhi-Jiang, HAY, Simon I. and DWYER-LINDGREN, Laura (2022). Mapping age- and sex-specific HIV prevalence in adults in sub-Saharan Africa, 2000–2018. BMC Medicine, 20 (1): 488. [Article]

Copyright and re-use policy

See http://shura.shu.ac.uk/information.html

¹ Additional File 1: Supplemental Information

- 2 Mapping age- and sex-specific HIV prevalence in adults in sub-
- 3 Saharan Africa, 2000–2018
- 4

5 Contents

6	1 Compliance with the Guidelines for Accurate and Transparent Health Estimates Reporting (GATHER)3
7	2 HIV data sources and data processing
8	2.1 Seroprevalence surveys
9	2.1.1 Data identification strategy6
10	2.1.2 Data processing for microdata6
11	2.1.3 Data processing for reports7
12	2.2 Antenatal care (ANC) sentinel surveillance7
13	2.2.1 Data sources7
14	2.2.2 Data processing
15	2.3 Polygon and age-aggregated data processing9
16	3 Covariate and auxiliary data10
17	3.1 Pre-existing covariates
18	3.2 Covariates constructed for this analysis11
19	3.2.1 Covariate selection criteria and definitions11
20	3.2.2 Covariate data13
21	3.2.3 Covariate modeling16
22	3.3 Administrative boundaries17
23	3.4 Gridded population17
24	4 Statistical model17
25	4.1 Covariate stacking17
26	4.2 Geostatistical model
27	4.2.1 Model description
28	4.2.2 Model fitting and prediction23
29	4.3 Model validation25
29 30	4.3 Model validation

32	4.3.3 Comparisons to adult prevalence estimates	29
33	4.4 Post-estimation	29
34	4.4.1 Aggregation to first- and second-level administrative subdivisions	29
35	4.4.2 Calibration to Global Burden of Disease 2019	
36	4.4.3 Calculating people living with HIV (PLHIV)	
37	5 References	

1 Compliance with the Guidelines for Accurate and Transparent Health

Estimates Reporting (GATHER)

Item #	Checklist item	Description of Compliance
Objectives and funding		
1	Define the indicator(s), populations (including age, sex,	Precision public health and HIV section
	time period(s) for which estimates were made.	
2	List the funding sources for the work.	Acknowledgments section
Data Inputs		
For all data inputs from multiple	sources that are synthesized as part	of the study:
3	Describe how the data were identified and how the data	Methods; Additional File 1: Sections 2.1, 2.2, 3.1, 3.2
	were accessed.	
4	Specify the inclusion and exclusion criteria. Identify all ad-	Additional File 1: Sections 2.1, 2.2, 3.1, 3.2, Additional File 2:
	hoc exclusions.	Table S3
5	Provide information on all	Additional File 2: Tables S1-2,4-
	included data sources and their	5,
	main characteristics. For each	https://ghdx.healthdata.org/rec
	data source used, report	ord/ihme-data/sub-saharan-
	reference information or	africa-hiv-prevalence-geospatial-
	contact name/institution,	estimates-2000-2018
	population represented, data	
	collection method, year(s) of	
	data collection, sex and age	
	range, diagnostic criteria or	
	measurement method, and	
	sample size, as relevant.	
6	Identify and describe any	Methods
	categories of input data that	
	have potentially important	
	blases (e.g., based on	
	characteristics listed in item 5).	
For data inputs that contribute to	o the analysis but were not synthesiz	ted as part of the study:
/	Describe and give sources for	Methods; Additional File 1:
	any other data inputs.	Sections 3.1, 3.3, 3.4;
For all data inputs:		
8	Provide all data inputs in a file	Available through
	format from which data can be	https://ghdx.healthdata.org/rec

	efficiently extracted (e.g., a	ord/ihme-data/sub-saharan-
	spreadsheet rather than a PDF)	africa-hiv-prevalence-geospatial-
	including all relevant meta-data	estimates_2000_2018
	listed in item 5. For any data	
	inputs that cannot be shared	
	hoose of othical or logal	
	because of ethical of legal	
	reasons, such as third-party	
	ownership, provide a contact	
	name or the name of the	
	institution that retains the right	
	to the data.	
Data analysis	1	
9	Provide a conceptual overview	Methods; Figure 2
	of the data analysis method. A	
	diagram may be helpful.	
10	Provide a detailed description of	Methods; Additional File 1:
	all steps of the analysis,	Sections 2-4
	including mathematical	
	formulae. This description	
	should cover, as relevant, data	
	cleaning, data pre-processing,	
	data adjustments and weighting	
	of data sources, and	
	mathematical or statistical	
	model(s).	
11	Describe how candidate models	Additional File 1: Section 4.3
	were evaluated and how the	
	final model(s) were selected.	
12	Provide the results of an	Additional File 1: Section 4.3
	evaluation of model	
	performance, if done, as well as	
	the results of any relevant	
	sensitivity analysis.	
13	Describe methods for calculating	Methods: Additional File 1:
	uncertainty of the estimates.	Sections 3.2. 4.4
	State which sources of	
	uncertainty were and were not	
	accounted for in the uncertainty	
	analysis	
14	State how analytic or statistical	Available through
¹ 7	source code used to gonorate	https://github.com/ibmouw/lbd
	estimates can be accessed	/tree/hiv_prev_africa_2020
Results and Discussion		

15	Provide published estimates in a	Available through
	file format from which data can	https://ghdx.healthdata.org/rec
	be efficiently extracted.	ord/ihme-data/sub-saharan-
		africa-hiv-prevalence-geospatial-
		estimates-2000-2018
16	Report a quantitative measure	Results; Figure 4; Additional File
	of the uncertainty of the	3: Figs. S27-34
	estimates (e.g. uncertainty	
	intervals).	
17	Interpret results in light of	Results; Methods
	existing evidence. If updating a	
	previous set of estimates,	
	describe the reasons for	
	changes in estimates.	
18	Discuss limitations of the	Discussion; Methods
	estimates. Include a discussion	
	of any modelling assumptions or	
	data limitations that affect	
	interpretation of the estimates.	

42 2 HIV data sources and data processing

43 2.1 Seroprevalence surveys

44 2.1.1 Data identification strategy

45 We identified HIV seroprevalence surveys in sub-Saharan Africa (SSA) through a review of all surveys in the Demographic and Health Survey (DHS), AIDS Indicator Survey (AIS), Multiple Indicator Cluster Survey 46 47 (MICS) series, and other surveys listed in the Global Health Data Exchange[1]; surveys included in the 48 national HIV estimates files from UNAIDS[2]; and surveys listed in the US Census Bureau HIV/AIDS 49 Surveillance Database[3]. For a survey to be considered for this analysis, we required that the survey 50 reported HIV blood test results, sampled from the general adult population, and contained geographic 51 information more refined than country level. For surveys with no microdata available we used reports if 52 they included sample size, or uncertainty intervals from which sample size could be derived. Our desired 53 age range was 15–59 years, but we also included survey reports that recorded prevalence for age spans 54 within that range. The surveys used in this analysis are listed in Additional File 2: Table S1 and visualized 55 in Figure 1. We additionally considered data sources identified through literature review; however, 56 because data from these sources predominantly did not match our inclusion criteria related to age 57 distribution (see section 2.1.3 below), we elected to exclude all literature review data from this model.

58 Other survey data exclusions are detailed in Additional File 2: Table S2.

59 2.1.2 Data processing for microdata

To prepare survey microdata for analysis, we first subset the data to the age range of interest, 15–59 60 years, and dropped any data that were not sex-specific. For data coded by gender rather than sex, we 61 62 treated these data as if they were sex-specific rather than gender-specific. We then dropped rows for 63 individuals explicitly listed as not tested or where the blood samples were marked as lost or rejected 64 (insufficient sample volume, tip broken, etc.). Inconclusive and indeterminate test results were coded as 65 a negative test result. After subsetting according to these conditions, we further dropped any microdata missing an HIV test result, survey weight, or geographic information or due to the GPS coordinates being 66 67 located more than 10 km outside of the country border. Coordinates within 10 km of the country border 68 were snapped to be approximately 1 km inside the nearest border of the specified country.

69 We then aggregated the individual-level microdata into sex-specific five-year age bins (15–19, 20–24,

70 25–29, 30–34, 35–39, 40–44, 45–49, 50–54, 55–59; hereafter termed 'ages') to the finest possible

representing the location of the survey

72 cluster (point-level data). The interview date for each specific location was calculated as the median of

73 the individual-level interview dates. Where point-level referencing was not available, we geolocated

survey microdata to the smallest geographical area (termed 'polygon') possible. Individual-level sample
weights were used when calculating prevalence, and the effective sample size for each prevalence
estimate was estimated via the Kish approximation[4], which accounts for differences in the underlying
selection probability within a sample.

78 2.1.3 Data processing for reports

79 In instances where individual-level microdata were not available, we used summary reports, given that 80 the estimates reported were similar in nature to what we would calculate from the microdata. We used 81 the median months of the reported data collection periods as the interview dates to align with the 82 extracted microdata. If sample sizes were not included in the report, we estimated them from the 83 reported confidence intervals, assuming that a Normal approximation was used to generate 95% 84 confidence intervals. In both instances, sample sizes were further adjusted by multiplying the median 85 design effect (ratio of effective sample size to observed sample size) calculated in the microdata as 86 described above. We only used reports with sex-specific estimates. Summary reports only provide 87 estimates aggregated across age; we included only those that completely covered either some or all the 5-year age bins within the 15–59 year age range being modeled. Because of their incongruity with our 88 89 methods for modeling age-aggregated data (detailed in Additional File 1: Section 2.3), we did not include 90 reports extending below 15 years or above 59 years, or any reports incompletely covering any of our 91 ages. For example, we included reports covering age ranges such as 15–59 years, or 15–49 years, but 92 excluded reports covering age ranges such as 15–64 years, or 18–24 years.

93 2.2 Antenatal care (ANC) sentinel surveillance

94 2.2.1 Data sources

In addition to general population surveys, we used antenatal care (ANC) sentinel surveillance data, which measure HIV prevalence among pregnant females attending antenatal care clinics. Most of these raw data came from national Spectrum files that were developed by a country team of experts and compiled and shared by the UNAIDS secretariat[2]. These files include the HIV prevalence and sample size of ANC sentinel surveillance and routine testing for various sites and years. We only used the sentinel surveillance estimates for our analysis.

- 101 We supplemented this data with ANC sentinel surveillance country reports. In general, the reports
- 102 contained the same information as the Spectrum files, but there was some additional information in the
- 103 reports and some discrepancies compared to the Spectrum files. The additional information included
- additional sites, additional years for given sites, and more precise prevalence estimates. In instances

where there were discrepancies for a given site-year, we elected to use the source where HIV
 prevalence was closest to the average prevalence of surrounding years for the same site.

107 Four countries had a notably large number of discrepancies between the Spectrum files and the ANC 108 reports. The Zambia Spectrum files recorded prevalence for the 15–39 years age range, while the 109 reports recorded prevalence for the 15–44 years age range. In this case, we elected to use the Spectrum 110 files because they had better data coverage in terms of number of site-years. There were also many 111 discrepancies in Central African Republic, Côte d'Ivoire, and Zimbabwe; we were unable to identify a 112 specific reason for these discrepancies and elected to use data from the Spectrum files only for these 113 countries. We investigated the ANC reports to determine if site names in the Spectrum files represented 114 hospitals, cities, or administrative subdivisions. We then used various mapping websites to find 115 geographic information related to these sites. For hospitals and cities/towns that are less than 25 km² in 116 area, we used a central GPS coordinate, and for administrative subdivisions we used a polygon of the 117 area. Some hospital sites had a city or town name rather than a hospital name. In those instances, we 118 searched for a hospital in the given city or town and used that hospital's GPS coordinates. If there were 119 multiple hospitals in the area but they were less than 5 km apart, we used the GPS coordinates of the 120 midpoint of the hospitals. If no hospitals were found in the area but the corresponding region was less 121 than 25 km² in area, we used the central GPS coordinate. Sites that could not be geolocated because 122 none of these conditions were met were excluded from further analyses.

123 2.2.2 Data processing

To prepare the ANC data for analysis, we compiled the HIV prevalence and sample size data from the Spectrum files and the ANC reports, and the site geographic information – either GPS coordinates or polygons for administrative subdivisions – into one dataset. After thoroughly inspecting the data, we decided to exclude the following data from our analysis:

- Hospital-level sites were dropped from Congo in 2011 (23 site-years) and Guinea-Bissau in 2003,
 2005, 2010, and 2014 (10 site-years) because the data aggregated by administrative subdivisions
 had better temporal coverage.
- We dropped administrative subdivisions that were masked by a different level of administrative
 subdivisions (8 site-years), defaulting to the level that would give better temporal coverage.
- We determined that 181 site-years were outliers based on inspection of site-level time trends
 and undue influence on model results, and these were dropped from the analysis.

135 Data from sites that could not be geolocated were also dropped (96 sites). Additionally, in the 136 Spectrum files, data from 12 site-years labeled as sentinel surveillance we suspect are actually 137 routine testing, and we excluded these from the analysis. In some cases, 'default' sample sizes 138 were reported for all sites and certain years in a given country (typically N = 300). In cases 139 where measured sample sizes were available for these affected sites in two or more other years, 140 we replaced the 'placeholder' sample size with the site-specific median from across measured 141 years. In cases where reported sample sizes were not available for other years, or where median 142 values clearly conflicted with site-specific trends in sample size over time, the 'placeholder' 143 sample size was retained. In the end, we adjusted sample sizes in this way for select data in five 144 countries, five years, and 57 sites, equating to 11 country-years and 146 site-years in total.

145 The ANC data included in this analysis are listed in Additional File 2: Table S2 and visualized in Figure 1.

146 2.3 Polygon and age-aggregated data processing

147 To incorporate observations geolocated to the polygon level as well as age-aggregated observations into 148 our model, we disaggregated these data to mimic point and/or age-specific data. Specifically, we 149 disaggregated each of these given observations to be location- and/or age-specific. For each polygon, 150 we generated points at the centroid of each pixel falling within that polygon and replicated that 151 observation's HIV prevalence and sample size at the location of each centroid. Age-aggregated data 152 were similarly disaggregated by replicating HIV prevalence and sample size once for each age covered in 153 the given age-aggregated observation's age range. In the cases of age-aggregated polygon data, these 154 two processes were combined. Next, each of the disaggregated, location- and age-specific rows of data 155 associated with a given aggregated observation were assigned weights (w_i) proportional to the age- and 156 sex-specific population at that location for the given year, derived from WorldPop[5]. For ANC data, ages 157 and locations within an ANC observation were weighted by births rather than population. The number 158 of births for a given age and location was calculated as the product of the location-, age-, and sex-159 specific population again derived from WorldPop[5], and the national fertility rate, derived from GBD 160 2019 estimates[6]. Weights per observation all summed to one. Age-specific point observations were 161 each assigned a weight of one.

To reduce the computational burden imposed by this method in terms of the large number of locations
and ages generated, in cases where for at least one location and/or age (*j*) within an observation,

164
$$w_j < \frac{1}{2} \cdot 1/\max(j),$$

- 165 we successively dropped the lowest-weighted locations and/or ages in that observation, until a
- 166 maximum of 1% of the observation's weight was dropped. Remaining locations and/or ages within that
- 167 observation were then reweighted to maintain a total observation weight of one. Age-specific point
- observations were each given a weight of one. This ultimately allowed us to retain ≥99% of our
- 169 observation weight of while removing 42.2% of pixel-ages, greatly mitigating the computational burden
- 170 on this model.

171 3 Covariate and auxiliary data

172 3.1 Pre-existing covariates

173 Mirroring the previously published adult HIV prevalence model[7], this analysis included five pre-existing 174 covariates: travel time to the nearest settlement of more than 50,000 inhabitants, total population, 175 night-time lights, urbanicity, and malaria incidence. These variables were selected from among available 176 gridded datasets for SSA because they are factors, or proxies for factors, that previous literature has 177 identified to be associated (not necessarily causally) with HIV prevalence. The first four variables were 178 included as measures or proxies for connectedness and urbanicity, as HIV historically spread through 179 SSA along travel routes[8, 9] and is typically found to be higher in more urban compared to more rural 180 locations. Malaria incidence was selected based on prior evidence relating higher malaria incidence 181 rates to higher prevalence of HIV at the population level[10, 11]. Sources for these data are given in 182 Additional File 2: Table S4. These covariates underwent spatial and temporal processing in preparation 183 for their inclusion in analysis. 184 Spatial processing involved resampling the input covariate raster to align the spatial resolution of the

covariate to the 5 x 5-km resolution used in modeling. For covariates that were originally at a finer
resolution, we resampled the raster by taking the neighborhood average (travel time to the nearest
settlement of more than 50,000 inhabitants, night-time lights, and urbanicity) or sum (total population)

- 188 of the finer covariate raster to produce one at a 5 x 5-km resolution. Malaria incidence was natively at a
- 189 5 x 5-km resolution and thus did not require additional spatial processing.

Temporal processing was required in instances where the original temporal resolution of the covariate was anything other than annual. To resolve from a coarser time period to an annual time period, we filled the intervening years with the value from the nearest neighboring year (urbanicity) or using an exponential growth rate model (total population). Night-time lights and malaria incidence were provided at a one-year temporal resolution and did not require interpolation. As travel time to the nearest settlement of more than 50,000 inhabitants was available only for a single representative year

- 196 (2015), this covariate was set to be unchanged over time. After interpolation, night-time lights and
- 197 urbanicity were still missing the most recent years of the 2000–2018 analysis period, and in these
- instances, we filled out the end of the time-series carrying forward the most recent year without
- 199 modification.
- **200** 3.2 Covariates constructed for this analysis

201 3.2.1 Covariate selection criteria and definitions

- 202 In addition to the five pre-existing covariates, we constructed eight additional covariates for this analysis 203 that were updated from the previously published adult HIV prevalence model[7]. Numerous studies 204 have been conducted in SSA on risk and protective factors for HIV infection, and these factors commonly 205 include sexual behavior and factors that are thought to influence the transmission of HIV during sexual 206 intercourse[12]. Potential covariates were informed by past literature and required to have a 207 demonstrated association with HIV prevalence, though not necessarily a causal relationship. 208 Furthermore, our selection of covariates depended on having adequate data coverage from data 209 sources that could be readily extracted. In total, eight covariates were constructed: 210 Prevalence of male circumcision, including medical or traditional circumcision ('male • 211 circumcision'); 212 Prevalence of self-reported STI symptoms (genital discharge and/or genital ulcer/sore) in the last 12 months ('STI symptoms'); 213 Prevalence of marriage or living with a partner as married ('in union'); 214 215 Prevalence of one's current partner living elsewhere among females ('partner away'); 216 Prevalence of condom use at last sexual encounter within the last 12 months ('condom last 217 time'); • Prevalence of sexual activity among young females ('had intercourse'); 218 219 Prevalence of males reporting multiple sexual partners within the last year ('multiple partners in 220 year'); 221 Prevalence of females reporting multiple sexual partners within the last year ('multiple partners') 222 in year').
- 223 The notion that male circumcision has a protective effect against acquiring HIV was first proposed in
- 1986, and since then more than 30 cross-sectional studies have found the prevalence of HIV to be
- significantly higher in uncircumcised males, as well as numerous prospective studies that have shown a
- protective effect ranging from 48% to 88%[13]. In 2005, following the interruption of a randomized,

227 controlled trial of male circumcision in South Africa that showed a 60% protective effect of circumcision, 228 WHO and UN agencies first acknowledged evidence of male circumcision's protective effect[14]. 229 Following these declarations, voluntary medical male circumcision clinics (VMMC) emerged as an HIV 230 prevention strategy in 15 countries in Eastern and Southern Africa with high HIV prevalence and low 231 levels of male circumcision[15]. Given male circumcision's linkage to HIV in the scientific literature, many 232 surveys record self-reported circumcision status. The modeling of male circumcision estimates in this 233 study closely mirrors the methods recently published by Cork *et al*[16]. Here, we extend the analysis to 234 include estimates for the year 2018, as well as additional countries included in this study but not in the 235 previous work.

236 Coinfection of HIV with viral and bacterial sexually transmitted infections (STIs), most notably herpes

237 simplex virus type 2, is a well-studied mechanistic factor associated with higher risk of HIV

acquisition[17]. STIs are thought to have been especially important risk factors during the early stages of

the epidemic when infections were concentrated in high-risk groups, though researchers have since

argued STIs are also critical in advanced stages[18]. Due to the association between STI prevalence,

241 sexual behavior, and HIV, most survey series detail the self-reported presence of STI symptoms,

facilitating its inclusion as an HIV covariate in this analysis.

Marital status represents a structural factor that, while distal to HIV exposure, has been associated with the number and type of sexual partners, as well as with HIV status[19, 20]. It has been postulated that the relationship between an individual's marital status and the number of sexual relationships regulates the protective effect of marriage on the risk of HIV infection[21]. Marital status is a readily available indicator in household surveys more generally.

248 The frequency with which a partner has slept away from home during the past year is an indicator of the

249 mobility of male partners, and studies have found that mobility confers an increased risk for HIV[22].

250 Part of the rapid spread of HIV in SSA has been attributed to occupations that consist of geographical

251 mobility, especially truck drivers, who are identified as high-risk for acquiring and spreading HIV[23].

252 Many surveys ask females if their partner has lived away from home in the past year, and we use these

253 responses as a proxy for occupational mobility.

254 Condom use is a sexual behavior factor that is protective against acquiring HIV. Condoms are often

presented as the most effective HIV prevention method of sexual transmission of the disease[24].

256 Though it is difficult to measure accurately how often condoms are used in sexual encounters, most

surveys report on the use of condoms in last sexual intercourse, a readily available proxy for overallcondom use.

An early age at sexual debut may be associated with the number of lifetime sexual partners, which is 259 260 considered a key risk factor for contracting HIV[21]. Furthermore, early age at sexual debut has been 261 shown to be associated with numerous other risk factors for HIV acquisition, such as STI prevalence and 262 decreased condom use[25]. For young females, the initiation of sexual activity is the first important 263 determinant of potential viral exposure, and delayed sexual debut has been associated with decreased 264 risk of HIV acquisition[26]. Given these relationships between HIV and age of sexual debut, and the 265 relative ease of acquiring self-reported sexual status, we constructed an indicator for whether young 266 (ages 15–24) females have had intercourse.

267 An individual's number of sexual partners correlates with HIV risk, and past studies have found a 268 relationship between the number of sexual partners and HIV prevalence[27]. The number of sexual 269 partners is thought to have been an especially important factor in the early stages of an epidemic, 270 though past research has determined it remains a key risk factor in advanced stages[18]. Surveys often 271 ask males and females their number of partners in the past year, and we used these responses to 272 construct a proxy for multiple concurrent sexual relationships. Separate covariates were constructed for 273 males and females given the well-documented discrepancy in the number of partners reported by males 274 as compared to females[28].

275 3.2.2 Covariate data

276 *3.2.2.1 Covariate data identification strategy*

277 We reviewed major survey series (Demographic and Health Surveys [DHS]; Multiple Indicator Cluster 278 Surveys [MICS]; AIDS Indicator Surveys [AIS]; Malaria Indicator Surveys [MIS]; Performance, Monitoring, 279 and Accountability Surveys [PMA]; Reproductive Health Surveys [RHS]; and Living Standards 280 Measurement Surveys [LSMS]) to identify surveys in SSA that contained relevant variables. We 281 supplemented this initial list of surveys with country-specific surveys identified in the Global Health Data 282 Exchange[1] and with a cross-check of all surveys extracted for HIV prevalence. We included surveys that 283 contain variables related to one or more of the covariate indicators (including any time restrictions 284 inherent to the indicator definition) and contained geographic information at a subnational level. 285 For all indicators except for 'had intercourse,' we required a survey to sample the general adult (ages 286 15–49) population. This age range was chosen for the covariates primarily due to data availability. For 287 'had intercourse', a survey only had to sample the general young female (ages 15-24) population to be

- included. Because covariates were not modeled to be age-specific, more discerning age range
- 289 requirements were not required of the surveys used in these models.
- 290 Because of variations we identified in the way these questions were asked across surveys, we tracked
- the skip logic and question format for all surveys including STI symptoms and/or the sexual activity
- 292 indicators. This helped us identify surveys for which the question format was so substantively different
- 293 from others as to require special handling or exclusion (e.g., questions asked without a time restriction
- for indicators that require a response from the last 12 months). We excluded select surveys because of
- these irreconcilable question variations, incomplete sampling (e.g., a specific age range or
- subpopulation), or untrustworthy or outlier data (as determined by the survey administrator or by
- inspection). The surveys used for these covariates are listed in Additional File 2: Table S5.

298 3.2.2.2 Covariate data processing for microdata

To prepare the survey microdata for analysis, we first constructed final indicators from the raw variablesincluded in the survey data:

- For 'STI symptoms,' we constructed a symptoms indicator that was true if a respondent
 reported either genital discharge or a genital sore/ulcer in the last 12 months, missing if either
 individual symptom was missing, and false if both symptoms were reported in the negative.
- For 'in union,' we constructed an indicator that was true for all respondents who reported being
 either currently married or living with a partner, false for any other marital status response, and
 missing if the marital status response was missing.
- For 'multiple partners in year', we used the reported number of sexual partners within the last
 12 months to construct a binary indicator that was true for any respondent reporting two or
 more partners and false for any respondent with 0 or 1 partners (including respondents who
 had never had intercourse).
- The other indicators were extracted from the survey microdata in their final form and required
 no additional construction.
- For each indicator, we subset the data to the desired age range (15–24 years for 'had intercourse', 15–
 49 years for all other indicators). For 'STI symptoms' we additionally restricted the sample to
 respondents who reported having had intercourse, while for 'partner away' we additionally restricted
 the sample to respondents currently 'in union'. We dropped any rows with missing responses or sample
 weights. For indicators where we model males and females together ('STI symptoms,' 'in union,'
- 318 'condom last time'), we dropped any surveys that did not interview both males and females. Any

observations missing geographic information or with inconsistent geographic information (i.e., points
 more than 10 km from the nearest specified country border) were also dropped.

321 Finally, we aggregated the weighted individual-level microdata for each indicator to the finest possible 322 spatial resolution available. We did not collapse or model covariate data according to specific age-bins 323 due to data limitations. As in Dwyer-Lindgren et al. [7], data for the covariate 'multiple partners per year,' was collapsed separately for males and females. 'Male circumcision' and 'prevalence of sexual 324 325 activity among young females' included data exclusively for males or females, respectively, but for all 326 other covariates, data were not collapsed to be sex-specific. Data were geolocated to latitude and 327 longitude at the survey cluster level wherever possible, and to the smallest possible polygon available 328 otherwise. As with the HIV prevalence data, we calculated the effective sample size for each spatial 329 aggregation using the Kish approximation[4].

330 *3.2.2.3 Covariate data processing for reports*

For 'male circumcision,' we also included summary reports for surveys where individual-level microdata were not available. We followed the same methods for report data processing as reported in Cork *et al*[16]. We chose not to include summary reports for other covariates. For 'STI symptoms,' the estimates included in reports used a different construction of the variable than that which we built from the microdata, making the reports incompatible with the microdata. For the sexual activity indicators, we decided against summary report extraction due to the significant number of surveys we were able to extract at the microdata level and the scarcity of reports for most of these indicators.

338 *3.2.2.4 Covariate data processing for polygons*

339 As with HIV prevalence data, wherever possible, covariate data were matched to a specific latitude and 340 longitude, and otherwise to the smallest areal unit (polygon) possible. The statistical model we 341 employed for covariate modeling required point-referenced data, so data matched to polygons were 342 resampled to generate pseudo-point data based on the underlying population distribution within the 343 polygon. The methods for the resampling are consistent with those previously used in the geospatial 344 modeling of many indicators, including adult HIV prevalence[7] and under-5 mortality[29]. Specifically, 345 for each polygon-level observation, we randomly sampled 10,000 locations among grid cells in the given 346 polygon with probability proportional to grid cell population. Grid cells were defined to be contained 347 within the polygon if their centroid fell within the geographic boundary. We performed k-means 348 clustering (with k set to 1 per 40 grid cells) on the sampled points to generate a reduced set of locations 349 to be used in modeling based on the k-means cluster centroids. Weights were assigned to each pseudo-350 point proportional to the number of sampled points contained in each of the k-means clusters, i.e., the

- number of sampled points divided by 10,000. Each pseudo-point generated by this process was assigned
- 352 the HIV prevalence observed for the polygon as a whole, and a sample size equal to the sample size for
- the polygon as a whole multiplied by the weight derived for each point.

354 3.2.3 Covariate modeling

- Each of these covariates was estimated using a simplified version of the modeling framework used for HIV prevalence as described in Additional File 1: Section 4.2, closely mirroring the framework previously used to model adult HIV prevalence[7]. Notable differences from the age- and sex-specific HIV prevalence model reported in this paper included:
- No covariates were included in the covariate geospatial models;
- No corrections for data derived from ANC sentinel surveillance were included (as no such data
 were used in these models);
- Covariate prevalence was modeled entirely at the disaggregated level (i.e., space- and time specific). This was possible for covariate models because prevalence was specified at the age aggregated level, and polygon data were resampled into pseudo-points;
- Because the covariate models did not include age or sex dimensions, only the spatiotemporal
 Gaussian process term was included;
- An unstructured error term (or 'nugget effect') for location *s* and year *t* was included;
- A fixed effect on time was included. This was particularly important for 'male circumcision' for
 capturing the growing emphasis on voluntary medical male circumcision as an intervention for
 HIV prevention[16]. For other covariates, this effect captured general regional time trends;
- Covariate models were fit in R-INLA[30]. Modeling in R-INLA was possible for the covariate
 models due to their more simplistic specifications relative to the age- and sex-specific HIV
 prevalence model.
- 374 Therefore, these models were specified as follows:
- 375 $Y_j \sim Binomial(N_j, p_j)$
- $logit(p_j) = \beta_0 + \beta_1 t + \gamma_{c[l]} + Z_j + \epsilon_j$
- 377 $\gamma_{c[l]} \sim Normal(0, \sigma 2_{country})$
- $Z_j \sim GP(0, \Sigma_{space} \otimes \Sigma_{time})$
- $\epsilon_j \sim Normal(0, \sigma_{nugget}^2)$

381	• N_j is the number of individuals sampled and Y_j is the number of individuals who tested positive,
382	or answered affirmatively among those sampled for the given covariate, for a given location and
383	year (j);
384	• p_j is the underlying prevalence for the given covariate for a given location and year j ;
385	• β_0 is an intercept;
386	• $\beta_1 t$ is a fixed effect for a given year t ;
387	• $\gamma_{c[l]}$ is a country-level random effect for country <i>c</i> containing location <i>l</i> ;
388	• Z_i is a spatially and temporally correlated random effect for a given location and year j ;
389	• ϵ_i is an independently distributed random effect for a given location and year <i>j</i> .
390	All priors and hyper-priors were otherwise the same as those used for the same respective terms in the
391	previously published adult HIV prevalence model[7]. Maps of each constructed covariate in 2000, 2005,
392	2010, and 2018 are displayed in Additional File 3: Figs. S1-8.
393	3.3 Administrative boundaries
394	For this analysis we used shape files from the Database of Global Administrative Areas (GADM)[31] to
395	define country boundaries and first- and second-level administrative subdivisions. We manually updated
396	known discrepancies.
397	3.4 Gridded population
398	The gridded population data used for this analysis were obtained from WorldPop[5]. Because WorldPop
399	provides data at a 1 x 1-km spatial resolution at five-year intervals, we processed these data as
400	described in Additional File 1: Section 3.1 to aggregate to a 5 x 5-km spatial resolution and interpolate to
401	annual time periods. When we use population as a covariate, we use total population. In all other
402	instances (as described in Additional File 1: Sections 2.3 and 4.4) we use age- and sex-specific
403	population.

404 4 Statistical model

380

where:

405 4.1 Covariate stacking

Stacked generalization/regression, or stacking, is an ensemble modeling method that combines multiple
 prediction methods to increase predictive validity relative to a single modeling approach. This ensemble
 modeling method relies on a variety of sub-models that are then combined by a secondary learner to

409 produce a meta-model that fuses multiple algorithmic methods to capture nonlinear effects and

410 complex interactions[32]. Our implementation of stacking largely follows the approach described by 411 Bhatt and colleagues[33] and which was previously implemented for modeling adult HIV prevalence[7]. 412 Because the HIV-specific covariates were modeled at the age- and (largely) sex-aggregated level, we fit 413 the stacker models at that same level, using HIV prevalence data aggregated across ages 15–49 years 414 and both sexes. The age range 15–49 years was used in this case because of its predominant use in 415 seroprevalence surveys compared to the 15–59 years range, allowing us to retain more data for use in 416 stacking purposes. Polygon data were excluded from stacking models due to their incongruity with the 417 configurations needed for the different sub-models. The ANC data were also excluded due to known 418 sampling biases, which are described in the Additional File 1: Section 4.2.

419 We fit three sub-models – a generalized additive model, boosted regression trees, and lasso regression –

420 to the HIV survey data with the five pre-existing and eight constructed covariates as well as calendar

421 year included as explanatory variables. We selected these three sub-models based on ease of

422 implementation through existing software packages, the fundamental differences in their approaches,

423 and a proven track record of predictive accuracy[33]. Sub-models were fit in R using the mgcv[34],

424 xgboost[35], glmnet[36], and caret[37] packages.

425 Each sub-model was fit using five-fold cross-validation to avoid overfitting, and hyper-parameter fitting 426 was done to maximize predictive power. For each sub-model, we produced two sets of predictions: out-427 of-sample and in-sample. Out-of-sample predictions for each model were generated by compiling the 428 predictions from the five holdouts from each cross-validation fold, and in-sample predictions were 429 generated by re-fitting the sub-models using all available data. The out-of-sample sub-model predictions 430 were used as explanatory covariates when fitting the geostatistical model described below, and the in-431 sample predictions were used when generating predictions from the geostatistical model in order to 432 maximize data use. In both instances, the logit-transformation of the predictions was used to put these 433 predictions on the same scale as the linear predictors in the geostatistical model. Maps of in-sample 434 predictions from each stacker are presented in Additional File 3: Figs. S9-11.

435 4.2 Geostatistical model

436 4.2.1 Model description

We modeled HIV prevalence using a generalized linear mixed effects model discretized by space, time,
age, and sex. To simultaneously model our point and polygon observations, and our age-specific and
age-aggregated observations, we modeled prevalence at the observation level (*i*). However, prevalence
was first specified at the space, time, age-, and sex-disaggregated level (*j*):

441

- 442 $Y_i \sim \text{Binomial}(N_i, p_i)$
- 443 $\operatorname{logit}(p_{j}) = \beta_{0} + \beta_{1}X_{j} + Z_{1,j} + Z_{2,j} + Z_{3,c[j]}$

444
$$Z_{1,j} \sim GP(0, \Sigma_{space} \otimes \Sigma_{1,time})$$

445
$$Z_{2,j} \sim \text{GMRF}(0, \Sigma_{2,time} \otimes \Sigma_{2,age} \otimes \Sigma_{2,sex})$$

446

447 where:

448	• N_i and Y_i are the number of individuals sampled and the number of individuals who are HIV ⁺
449	among those sampled, respectively, at the observation level (i) ;
450	 <i>p_i</i> is the underlying HIV prevalence at the observation level <i>i</i>;
451	• p_j is the underlying HIV prevalence at the fully disaggregated (i.e., location, year, age, and sex-
452	specific; j) level;
453	• β_0 is an intercept;
454	• X_j is a vector of logit-transformed stacked covariates at the disaggregated level j , and $oldsymbol{eta}_1$ is the
455	corresponding vector of regression coefficients;
456	 Z_{1,j} random effects correlated across space and time;
457	• $Z_{2,j}$ is a random effect correlated across time, age, and sex;
458	• $Z_{3,c[j]}$ is a country-specific (<i>c</i>) random effect correlated across age.
459	Descriptively, we modeled the number of HIV-positive individuals (Y_i) among a sample (N_i) for a given
460	observation i as a binomial variable. The model first specified logit-transformed prevalence at the
461	disaggregated level (p_j) as a linear combination of a regional intercept (eta_0), age- and sex-specific
462	covariate effects ($oldsymbol{eta}_1 X_j$), and random effects correlated across space, time, age, and sex

 $Z_{3,c[i]} \sim \text{GMRF}(0, \Sigma_{3,c})$

- 463 $(Z_{1,i}, Z_{2,i}, Z_{3,c[i]})$. The intercept captures the overall mean level of HIV prevalence, while the covariate
- 464 effects capture the spatial and temporal variation in HIV prevalence that can be described as a function
- 465 of spatial and temporal variation in the included covariates. The random effects correlated across space,
- 466 time, age, and sex capture additional variation by location (within and between countries), time, age,
- 467 and sex that varies smoothly over these dimensions.
- 468 We then applied age-specific transformations related to fertility to p_i (described below), calculated as:

469
$$p_{transformed,j} = \frac{(p_j \cdot FRR_j)}{(p_j \cdot FRR_j) + 1 - p_j}$$

470 where:

471	٠	$p_{transformed,j}$ is the underlying HIV prevalence at the disaggregated level j , transformed to
472		account for age-specific differences in fertility within observation-level data derived from
473		antenatal care clinic sentinel surveillance. For all other survey data, $p_{transformed,j}$ = p_j ;
474	•	And FRR_j is the fertility rate ratio between HIV ⁺ and HIV ⁻ females at the disaggregated level j ,
475		used to correct for age-specific differences within observation-level (i.e., in this case, age-
476		aggregated) data derived from data derived from antenatal care clinic sentinel surveillance. For
477		all other survey data, $FRR_j = 1$;
478	Fir	ally, prevalence at the observation level (p_i) was then specified as:

479
$$p_{i} = \text{logit}^{-1} \left(\text{logit} \left(\sum \left(p_{transformed,j} \cdot w_{j} \right) \right) + \left(\beta_{2} + U_{s[i]} \right) \cdot I_{ANC} + \epsilon_{i} \right)$$

480
$$U_{s[i]} \sim \operatorname{Normal}(0, \sigma_{site}^2)$$

481
$$\epsilon_i \sim \operatorname{Normal}(0, \sigma_i^2)$$

482 where:

- 483 w_j is the weight applied to data at the disaggregated level. For point and age-specific data, $w_j =$ 484 1;
- I_{ANC} is an indicator variable that is 1 for data derived from antenatal care clinic sentinel
 surveillance and 0 otherwise;
- 487 β_2 is a fixed offset for observation-level data derived from antenatal care clinic sentinel 488 surveillance;
- 489 U_{i[s]} is a site-level random effect for data derived from antenatal care clinic sentinel surveillance
 490 for observation *i* containing ANC site *s*;
- 491 and (ϵ_i) is an observation-level error term.
- 492 Technically our polygon and age-aggregated data would follow a convolution of a mixture of binomial
- 493 distributions. However, for computational efficiency we instead implement here a binomial

494 approximation where for a given observation *i*:

$Y_i \sim \text{Binomial}(N_i, p_i)$

$$p_i = \sum_j w_j p(x_j) / \sum_j w_j$$

495 where we take w_i to be the population density proportion at pixel-age j (i.e., location and age x_i) for the polygon and/or age range for observation *i*, and $\sum_{i} w_{i} = 1$. We expected increased variance in our 496 497 estimates given this modeling framework compared to a model with equal data coverage that used only 498 point and age-specific data; however, given the limited availability of point and age-specific data, 499 sensitivity analyses (see Additional File 1: Section 4.3 and Additional File 3: Figs. S13-15) demonstrate 500 the larger benefit to our model in terms of reducing bias and error provided by the inclusion of 501 aggregated data. We chose this method for including aggregated data rather than the polygon 502 resampling method previously used to model adult HIV prevalence[7] among other indicators because 503 polygon resampling is less robust[38], isn't able to account for variation in the spatial covariates or 504 spatial field within polygon data sources, and uses an ad-hoc method for down-weighting the sample 505 size of the resampled points. Also, the new method enabled us to disaggregate data not only over space 506 but also by age, and allowed us to account for ANC-related bias at both age-aggregated and age-507 disaggregated levels.

508 HIV prevalence as measured by sentinel surveillance of antenatal care (ANC) clinics is known to be 509 biased as a measure of HIV prevalence in the general adult population because it captures pregnant 510 females who attend ANC only, as compared to all adult females [39, 40]. This bias may be either positive 511 or negative: the fact that all pregnant females are sexually active tends to elevate their risk of having 512 acquired HIV prevalence compared to the general female population (some of whom are not sexually 513 active), while HIV-related sub-fertility tends to reduce the prevalence of HIV⁺ females among the 514 population of pregnant females[41, 42]. Additionally, HIV-related sub-fertility tends to vary across 515 ages[43]; however, ANC data reported at the age-aggregated level does not account for these 516 differences. Further, we do not expect the sampling bias within age- and spatially aggregated ANC 517 observations to correspond with underlying populations, as we do for survey data. Nevertheless, ANC 518 data have better temporal and spatial coverage in many countries than survey data alone (Figure 1). We 519 therefore incorporated ANC data to capitalize on this additional data coverage, but also attempted to 520 correct for the known biases in multiple ways.

First, to account for age-specific differences in the fertility rate ratio of HIV⁺ and HIV⁻ females, we
corrected prevalence estimates from ANC clinics at the disaggregated level according to age-specific
fertility rate ratios, calculated according to age-specific and HIV-status-specific fertility estimates from

524 GBD 2019[6]. Fertility rate ratios were calculated at the national level, except for in Ethiopia, Kenya,

525 Nigeria, and South Africa, where estimates were available at the first administrative level.

526 Second, because we expect sampling prevalence for ANC data disaggregated over space and age to vary

527 according to age- and location-specific ANC clinic visitation rates, rather than according to the

528 distribution of the underlying population, we calculated the w_i values for disaggregated ANC data to

reflect this. Specifically, we used the number of births in a given year, location, and age as a proxy for

530 ANC visitation rate, and weighted disaggregated ANC data accordingly. Births were calculated by

531 multiplying the local population of females in the given year and age (based on local estimates from

WorldPop[5]) by the national fertility rate for that year and age (based on national-level estimates fromGBD 2019[6]).

534 Third, we accounted for ANC-related bias at the observation level. In instances where data in our model 535 were derived from ANC sentinel surveillance (I_{ANC} = 1), our model allows for this bias via a fixed term 536 (β_2) that captures the overall mean bias, and a site-specific random effect $(U_{i[l]})$ that captures local 537 differences in the extent of this bias. This approach is conceptually like previously described approaches 538 for spatial modeling using non-randomized (and therefore potentially biased) data and randomized 539 survey data[44, 45]. Although the bias associated with ANC sentinel surveillance may also vary over time 540 in addition to varying spatially, we felt there was insufficient data to estimate both spatial and temporal 541 variation in this bias, and so the bias associated with ANC sentinel surveillance was assumed to be time-542 invariant over the period of this analysis.

The spatially and temporally correlated random effect (Z_{1j}) was modeled as a Gaussian process with mean 0 and a covariance matrix given by the Kronecker product of a spatial Matérn covariance function (Σ_{space}) and a temporal first-order autoregressive (AR1) covariance function (Σ_{time}) . The Matérn covariance function is given by:

547
$$\Sigma_{space} = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \cdot (\kappa D)^{\nu} \cdot K_{\nu}(\kappa D)$$

In this analysis v (the smoothness parameter) was fixed at 1. A penalized complexity (PC) prior was used for the Matérn covariance function and specified via two hyper-parameters: the spatial range, ρ_s (where $\rho_s = \sqrt{8\nu}/\kappa$ and is equal to the distance at which correlation is approximately 0.1; the subscript *s* for space is used as to not confuse with the other correlation parameters, below), and marginal standard deviation, σ . PC priors shrink towards a more simplistic base model – in this case, one where the marginal variance is 0 and the spatial range is infinite – and are specified via setting the tail probabilities on each hyper-parameter[46, 47]. We followed the guidance provided by Fugulstad et al., who recommend selecting priors that satisfy $P(\sigma > \sigma_0) = 0.05$ and $P(\rho_s > \rho_{s_0}) = 0.05$, where σ_0 is between 2.5 to 40 times the expected true marginal standard deviation and ρ_{s_0} is between 1/10 to 1/2.5 of the expected true range[48]. Specifically, we set:

558 $\sigma_0 = 5; P(\sigma > \sigma_0) = 0.05$

559
$$\rho_{s_0} = 0.01 \ radians; \ P(\rho_s > \rho_{s_0}) = 0.05$$

560 Separate σ parameters were specified for each Z_j term included in the model; each was assigned the 561 same prior as above. Individual σ parameters were also included for the observation-level error term

562 (ϵ_i) and the ANC random effect ($U_{s[i]}$), with respective priors set as:

563
$$\sigma = 3; P(\sigma > \sigma_0) = 0.05$$

Additionally, for all Z_j terms included in the model, the AR1 covariance function is associated with different parameters accounting for correlations in time, age, and sex $-\rho_t$, ρ_z , and ρ_x , respectively. Unique ρ parameters were identified in each of their respective appearances in the model. For example, because an AR1 temporal covariance function was incorporated into the covariance matrices for $Z_{1,j}$, and $Z_{2,j}$, we fit two separate ρ_t parameters ($\rho_{1,2,t}$). We nevertheless used the same following hyperprior for all ρ_t , ρ_z , and ρ_x parameters, which corresponds to a prior mean of 0.76 with a 95% range of -0.17 to 0.97:

571
$$\log\left(\frac{1+\rho}{1-\rho}\right) \sim \operatorname{Normal}(2, 1.2^2)$$

- 572 Finally, priors for fixed effects were set as:
- 573 $\beta_0 \sim \text{Normal}(0, 3^2)$

574
$$\beta_1 \sim \text{Normal}(0, 3^2)$$

575
$$\beta_2 \sim \operatorname{Normal}(0, 3^2)$$

576 4.2.2 Model fitting and prediction

577 This model was fit in Template Model Builder (TMB)[49], package in R version 3.6.1. We used the

578 stochastic partial differential equations (SPDE) approach[50] to approximate the continuous

- 579 spatiotemporal Gaussian random field $(Z_{1,i})$. We constructed a finite elements mesh for the SPDE
- 580 approximation to the Gaussian process regression using a simplified polygon boundary (Additional File 3:

Fig. S41). We used a spatial mesh that was constructed on the S² domain which allowed distance to be
calculated along a sphere instead of using Euclidean distance between latitude and longitude
coordinates. We set the inner mesh triangle minimum edge length to 35 km, the maximum triangle
length to 500 km, with the mesh extending 500 km past the region's boundary. We used maximum a
posteriori (MAP) inference, using a maximum likelihood estimation with an augmented optimization
objective (log-likelihood function) which incorporated prior distributions for all model parameters.
Estimated model parameters are listed in Additional File 2: Table S6.

588 Due to computational constraints and to allow for regional differences in the relationship between 589 covariates and HIV prevalence as well as the strength of auto-correlation across space, time, age, and 590 sex in HIV prevalence, separate models were fit for four regions (Additional File 3: Fig. S12). Specifically, 591 we used the regional classifications for SSA from the Global Burden of Disease (GBD) study[51] which 592 group countries by location and epidemiological profile. We made small modifications to this 593 classification, grouping Sudan as part of the Eastern SSA region rather than the North Africa and the 594 Middle East region. We also dropped Cape Verde, Comoros, São Tomé and Príncipe, and Mauritania 595 from these modeling regions due to data missingness.

After fitting each model, we generated 1,000 draws of all model parameters from the approximated
joint posterior distribution using a multivariate-normal approximation. For each draw *s* of the model
parameters, we constructed a draw of

599
$$p_j^{(s)} = \log i t^{-1} \left(\beta_0^{(s)} + \beta_1^{(s)} X_j + Z_{1,j}^{(s)} + Z_{2,j}^{(s)} + Z_{3,c[j]}^{(s)} \right)$$

600 I_{ANC} is set to 0 for the purposes of generating estimates, so draws of β_2 and U_i are not incorporated 601 when generating draws of p_j . Additional processing of the output from the multivariate-normal 602 approximation is required for the spatial-temporal random effect $(Z_{1,j}^{(s)})$ prior to constructing 603 $p_j^{(s)}$ according to the equation above. Specifically, for $Z_{1,j}^{(s)}$, draws are generated initially only at vertices 604 of the finite element mesh, so we project from this mesh to each pixel-year combination desired for 605 prediction, i.e., the centroid of each grid cell on a 5 × 5-km grid as well as all years from 2000 to 2018. At 606 the end of this process, we have 1,000 draws of $p_j^{(s)}$ for each grid cell and year combination.

607 4.3 Model validation

628

608 4.3.1 Validation strategy

609 We used five-fold out-of-sample cross-validation in order to assess the performance of the modeling 610 framework described above with respect to predicting HIV prevalence. We first split all location- and 611 age-specific data into five groups using spatial and temporal stratification[52]. Temporal folds were 612 created by stratifying across years such that each fold contains approximately 1/5 of the data for each 613 year. Spatial folds were constructed we used a modified quadtree algorithm to spatially aggregate data 614 points. This algorithm recursively partitions two-dimensional space, alternating between horizontal and vertical splits on the weighted data sample size medians. The depth of recursive partitioning is 615 616 constrained by the target sample size within a partition and the minimum number of clusters or pseudo-617 clusters allowed within each spatial partition. The minimum sample size was set according to data 618 availability in each region—the minimum sample size was set at 425 for Central SSA and Southern SSA, 619 and 500 for Eastern SSA and Western SSA. These partitions were then allocated to one of five folds for 620 cross validation. This resulted in five groups that are approximately equal in terms of the total effective 621 sample size. We then fit the model described above five times, excluding each of the five holdout data 622 groups in turn. All ANC data were included in all models and were not used to assess model 623 performance given the known biases in these data. Due to difficulties in comparing age-aggregated and 624 polygon data to age- and location-specific results, polygon and age-aggregated survey data were 625 excluded from use in assessing model performance and were therefore used in all models. 626 After fitting the model five times, the data withheld from each model were matched with predictions from that model, and then these data-prediction pairs were compiled across all five models, resulting in 627

629 included in the analysis. HIV prevalence estimates based on single survey clusters are generally quite

a complete dataset of out-of-sample predictions corresponding to all location- and age-specific data

noisy due to very small sample sizes and are consequently insufficient as a 'gold standard' for evaluating

the model predictions[29]. To address this issue, we aggregated both the observed data and the

632 corresponding age- and sex-specific out-of-sample predictions within countries and within first- and

633 second-level administrative subdivisions, by calculating a weighted mean of each using the effective

634 sample sizes as the weights. Then, across all data-estimate pairs, we calculated two summary measures:

the mean error (ME, a measure of bias) and the root-mean square error (RMSE, a measure of totalvariance).

In addition, for each data-estimate pair, we constructed 95% prediction intervals from the 2.5th and
97.5th percentiles of 1,000 draws from a binomial distribution corresponding to each of the 1,000

639 posterior draws of HIV prevalence with p equal to HIV prevalence in a given posterior draw and N equal

- to the effective sample size for the data point. We then calculated coverage as the percentage of data-
- 641 estimate pairs where the data point was contained within this 95% prediction interval. Finally, to
- 642 complement the out-of-sample predictive validity metrics, we calculated in-sample predictive validity
- 643 metrics using the same process but matching each data point to predictions from a model fit using all
- 644 data.

645 4.3.2 Sensitivity analyses

We used this validation strategy to assess model performance of the final model compared models of
adult prevalence, as well as a number of alternatives related to data inclusion and model
specification[53].

649 4.3.2.1 Adult prevalence sensitivity

650 We assessed the performance of our age- and sex-specific model compared to an adult-level HIV 651 prevalence model, that is, one for combined sexes and ages 15–49 years. In these comparisons, we 652 validated the results of the age and sex model not only at the age- and sex-disaggregated level, but also 653 for estimates re-aggregated to the adult level (see Additional File 1: Section 4.3.3). The adult prevalence 654 model we tested mirrored the age- and sex-specific model as closely as possible; all survey microdata 655 and reports for ages 15–49 years were included, as well as all ANC data. All parameters from the age-656 and sex-specific model were retained in the adult prevalence model, except those that pertained to age 657 and sex correlations (i.e., $Z_{2,j}$ and $Z_{3,[c]j}$). To replace the country-level variation provided in the age- and 658 sex-specific model by the country-specific age correlation term $(Z_{3,c[j]})$, we instead included a countrylevel random effect, $\gamma_{[c]j}$. Logit-transformed disaggregated prevalence $\mathrm{logit}(p_j)$ was therefore specified 659 660 as:

661

 $logit(p_i) = \beta_0 + \beta_1 X_i + Z_{1,i} + \gamma_{[c]i}$

663 Observation-level adult prevalence (p_i) was calculated using the same equation from age- and sex-664 specific prevalence estimation, differing only in that the transformation related to age-specific fertility-665 rate ratios (*FRR*) was not applied.

To assess our decision to employ novel methods for including polygon data in our model rather than the
 previously utilized polygon resampling technique[7, 54], we also compare our results to those of an
 adult prevalence model built using polygon resampling. We elected to test polygon resampling in an

age-aggregated model due to the age- and sex-specific model's heavy reliance on age-aggregated data,
which is processed in effectively the same manner as the polygon data. We therefore avoid this conflict
by testing resampling in the adult prevalence model. In total this resulted in the comparison of four
models and corresponding sets of results:

673 1. The final age- and sex-specific model, with age- and sex-specific results;

2. The final age- and sex-specific model, results re-aggregated to the adult level;

675 3. Results for an adult prevalence model, employing the novel polygon processing system as in the676 final model;

677 4. Results for an adult prevalence model, employing the previously published polygon resampling678 system.

679 Comparisons of adult prevalence when modeled versus re-aggregated can be seen in Additional File 3:

Fig. S16. The results of this sensitivity analyses can be found in Additional File 3: Fig. S13. The re-

aggregated adult estimates were outperformed by the modeled adult estimates in some respects, but

not in others. For example, our mean error calculations were much closer to zero (indicating less bias)

683 for modeled adult prevalence compared to re-aggregated estimates. This may ultimately be a product of

684 our process for re-aggregating age- and sex-specific estimates—these calculations are heavily influenced

by local population structure. We also calculated consistent overestimations for 95% coverage for the

age- and sex-specific model, indicating some overestimation of our uncertainty intervals compared to

687 modeled adult prevalence. Meanwhile RMSE tended to be substantially lower for the re-aggregated

688 estimates (indicating lower variance). The in- vs. out-of-sample results also tended to be more similar

689 within the re-aggregated estimates compared to other models, although this also varied by region.

690 Some necessary differences in data and model configuration likely contributed to these differences.

691 Further investigation of the influences on these differences will be an important future direction in this

692 line of research.

693 4.3.2.2 Data sensitivity

To assess the contribution of our different data sources, we tested additional models with the followingsubsets of the data:

696 1. Survey data only (no ANC data);

697 2. Point and age-specific data only (no polygon or age-aggregated data).

698 The results of this sensitivity analyses can be found in Additional File 3: Fig. S14. We found the 699 performance of these models using smaller data subsets to be very region-specific. For example, when 700 ANC data were excluded, mean error in Eastern and Southern SSA tended to be closer to zero (i.e., less 701 biased) compared to when these data were included. When all polygon and age-aggregated data were 702 excluded, Eastern SSA was still less biased, but in this case out-of-sample Southern SSA performed worse 703 than when all data were included. Central and Western SSA, on the other hand, performed dramatically 704 worse in terms of mean error when ANC as well as all polygon and age-aggregated data were excluded. 705 Survey data were severely limited in Central SSA in particular, so it is not surprising that estimates in this 706 region were highly dependent on ANC data. These results were similar for our other validation metrics. 707 Given that Eastern and Southern SSA have relatively better spatial and temporal survey data coverage 708 (Figure 1), it is expected that these regions would be more robust to the loss of ANC and other 709 aggregated data. It is clear that while these unconventional data sources provide tremendous insight in 710 the absence of better survey data coverage, more work is needed to reduce bias associated with their 711 inclusion.

712 4.3.2.3 Statistical configuration sensitivity

To assess our final chosen statistical configuration, we assess the utility of each term included in the

model by testing models excluding individual parameters. This resulted in six additional models:

715 1. No interaction between the space and time correlation terms;

716 2. No interaction between the time, age, and sex correlation terms;

3. No interactions whatsoever between the space, time, age, and sex correlation terms;

718 4. No country-specific age correlation term;

719 5. No observation-level error term;

720 6. No stackers.

721 In cases where interactions between terms were removed, the individual terms were retained if not 722 included elsewhere in the model. For example, for the model where the interaction between space and 723 time correlations was removed, and additional "space-only" correlation term was included, but because 724 the time correlation was still accounted for in the time-age-sex interaction, no additional time 725 correlation was included. The results of this sensitivity analyses can be found in Additional File 3: Fig. 726 S15. We note that in a number of respects, our final chosen model did not out-perform those excluding 727 some of our chosen parameters (Additional File 3: Fig. S15). For example, the out-of-sample RSME 728 values for our final model in many cases were higher than those for other tested models. We believe

729 this may be partially driven by the fact that our validation analyses are conducted exclusively for our 730 point and age-specific data. Given the heavy reliance of this model on polygon and age-aggregated data, 731 we believe that these sensitivity analyses provide an incomplete assessment of our model performance. 732 In in-sample testing, our final model did outperform other models with regards to RMSE, though we 733 acknowledge this does not speak to our ability to predict to sparsely sampled location-years. We also 734 found that in inclusion of some terms, such as the country-specific age effect, $Z_{3,c[i]}$ and the 735 observation-level error term helped to reduce bias and smooth trends at the national level, which may 736 not be reflected in these validation metrics. With additional data and computing power, it is probable 737 that this model would benefit from additional and more complex interactions. However, given the 738 resources currently available to us, we are confident that our final model represents the best possible 739 option at this time.

740 4.3.3 Comparisons to adult prevalence estimates

741 As this age- and sex-specific HIV prevalence model serves as a follow-up to a previously described 742 analysis of adult (ages 15–49 years) HIV prevalence[7], it was important that we compare the estimates 743 from this model to one mirroring its predecessor. To make this comparison effectively, it was necessary 744 that we re-aggregate our age- and sex-specific results to the 'adult' level. We therefore calculated HIV 745 prevalence for adults ages 15–49 years by summing our final age- and sex-specific PLHIV estimates 746 across males and females age groups 15–49 years, for each grid cell and year, and dividing those by cell-747 and year-specific population estimates summed across the same age groups. Both PLHIV and population 748 estimates were derived during the post-estimation process, described below in Additional File 1: Section 749 4.4. In select grid cells where the population was estimated to be zero, prevalence was weighted by the 750 second administrative-level population age and sex structure. For a description of the calculation of 751 second administrative-level estimates, see Additional File 1: Section 4.4. For sensitivity analyses, re-752 aggregated estimates were compared to location-specific survey microdata collapsed across all adults 753 ages 15–49 years, the same data used to validate modeled adult prevalence. For a comparison of these 754 results re-aggregated across sexes and age groups to HIV prevalence estimates modeled across adults, 755 see Additional File 3: Figs. S13 and S16.

756 4.4 Post-estimation

757 4.4.1 Aggregation to first- and second-level administrative subdivisions

758 In addition to estimates of HIV prevalence on a grid, we also constructed estimates of HIV prevalence for

- 759 first- and second-level administrative subdivisions. These estimates were derived by calculating
- population-weighted averages of HIV prevalence for each grid cell or fractional grid cell within a given

761 first- or second-level administrative subdivision for a given age, sex, and year. Grid cell fractions were 762 assigned at the second-level administrative subdivision shape to determine what fraction of the area of 763 each grid cell fell within each administrative unit. Since all second-level subdivisions nest within first-764 level subdivisions, which in turn nest within countries, this strategy assigned the cell fractions to an 765 administrative area at each level of the administrative hierarchy. We assumed that population density 766 within each cell was uniform, and for cells that were split across multiple subdivisions, allocated the 767 WorldPop population estimate in proportion to area. This process was carried out separately for each 768 modeling region, so cells that cross international borders that are also regional borders were allocated 769 in their entirety to the country that contained the centroid of the grid cell. This was carried out for each 770 of the 1,000 posterior draws at the grid cell level, generating 1,000 posterior draws for each 771 administrative subdivision. Final estimates and uncertainty intervals for each subdivision at each level of 772 the administrative hierarchy were derived from the mean, 2.5th percentile, and 97.5th percentile of 773 these draws, respectively.

4.4.2 Calibration to Global Burden of Disease 2019

775 To take advantage of the more epidemiologically structured modeling approach and additional national-776 level data used by GBD 2019, we performed post-hoc calibration of our estimates to the GBD 777 estimates[43]. Using the assignment of cells and cell fractions to the administrative hierarchy described 778 above, we first scaled the grid cell-level WorldPop estimates[5] to match the corresponding GBD 779 population estimates[6] for each country, year, age, and sex. To do so, for each country, year, age, and 780 sex, we defined a population raking factor as the ratio of the GBD population estimate to the sum of the 781 WorldPop population estimates for all cells and fractional cells within the country, and then multiplied 782 the WorldPop population estimates for all cells and fractional cells within the country by this raking 783 factor.

784 We then similarly adjusted our HIV prevalence estimates. Specifically, for each country, year, age, and 785 sex, we defined a prevalence 'raking factor' as the ratio of the GBD prevalence estimate to the 786 population-weighted mean of estimates for all cells and fractional cells within the country, and then 787 multiplied each HIV prevalence draw for all cells and fractional cells within the country by this raking 788 factor. At this point, the prevalence estimates for cells that had been fractionally allocated to multiple 789 countries were recombined by calculating a weighted average, with weights determined by the relative 790 area of each fraction. Final calibrated estimates for each grid cell were calculated as the mean of the 791 scaled draws, and 95% uncertainty intervals were calculated as the 2.5th and 97.5th percentiles of the

- scaled draws. The impact of this calibration procedure is depicted in Additional File 3: Figs. S17 and S18,
- which compares the pre-calibration estimates to the post-calibration estimates.
- **794** 4.4.3 Calculating people living with HIV (PLHIV)
- 795 We estimated the number of people living with HIV (PLHIV) in each grid cell, year, age and sex by
- 796 combining estimated population and HIV prevalence after calibration to GBD 2019 estimates as
- described above. Specifically, for each cell and fractional cell, we multiplied the estimated population by
- each of the 1,000 prevalence draws to generate 1,000 draws of PLHIV. Fractional cells were then
- recombined by summing PLHIV for each draw within each cell. Final point estimates and uncertainty
- 800 intervals for PLHIV were calculated as the mean, 2.5th percentile, and 97.5th percentile of these draws,
- 801 respectively.

802 5 References

- 1. Global Health Data Exchange | GHDx. Available at: http://ghdx.healthdata.org/. (Accessed: 16th June
 2020)
- 805 2. UNAIDS. National HIV estimates file. 2019.
- 806 https://www.unaids.org/en/dataanalysis/datatools/spectrum-epp.
- 3. The United States Census Bureau. HIV/AIDS Database. 2019. https://www.census.gov/programs surveys/international-programs/about/hiv.html. Accessed 16 Jun 2020.
- 4. Wiegand H. Kish, L.: Survey Sampling. John Wiley & Sons, Inc., New York, London 1965, IX + 643 S., 31
 Abb., 56 Tab., Preis 83 s. Biom Z. 1968;10:88–9.
- 5. WorldPop. WorldPop Dataset. 2020. http://www.worldpop.org.uk/data/get_data/.
- 812 6. GBD 2019 Demographics Collaborators. Global age-sex-specific fertility, mortality, healthy life
- 813 expectancy (HALE), and population estimates in 204 countries and territories, 1950–2019: a
- comprehensive demographic analysis for the Global Burden of Disease Study 2019. Lancet.
- 815 2020;396:1160–203.
- 7. Dwyer-Lindgren L, Cork MA, Sligar A, Steuben KM, Wilson KF, Provost NR, et al. Mapping HIV
 prevalence in sub-Saharan Africa between 2000 and 2017. Nature. 2019;570:189–93.
- 818 8. Tatem AJ, Hemelaar J, Gray RR, Salemi M. Spatial accessibility and the spread of HIV-1 subtypes and 819 recombinants. AIDS. 2012;26:2351–60.
- 9. Gray RR, Tatem AJ, Lamers S, Hou W, Laeyendecker O, Serwadda D, et al. Spatial phylodynamics of
 HIV-1 epidemic emergence in east Africa. AIDS Lond Engl. 2009;23:F9–17.
- Abu-Raddad LJ, Patnaik P, Kublin JG. Dual Infection with HIV and Malaria Fuels the Spread of Both
 Diseases in Sub-Saharan Africa. Science. 2006;314:1603–6.

- 11. Cuadros DF, Branscum AJ, Crowley PH. HIV–malaria co-infection: effects of malaria on the
 prevalence of HIV in East sub-Saharan Africa. Int J Epidemiol. 2011;40:931–9.
- 12. Auvert B, Buvé A, Ferry B, Caraël M, Morison L, Lagarde E, et al. Ecological and individual level
 analysis of risk factors for HIV infection in four urban populations in sub-Saharan Africa with different
 levels of HIV infection. AIDS. 2001;15:S15.
- 13. Bailey RC, Moses S, Parker CB, Agot K, Maclean I, Krieger JN, et al. Male circumcision for HIV
 prevention in young men in Kisumu, Kenya: a randomised controlled trial. Lancet Lond Engl.
 2007;369:643–56.
- 14. World Health Organization. WHO | UNAIDS statement on South African trial findings regarding male
 circumcision and HIV. 2005. https://www.who.int/mediacentre/news/releases/2005/pr32/en/.
- 15. Sgaier SK, Reed JB, Thomas A, Njeuhmeli E. Achieving the HIV prevention impact of voluntary
- medical male circumcision: lessons and challenges for managing programs. PLOS Med.
- 836 2014;11:e1001641.
- 16. Cork MA, Wilson KF, Perkins S, Collison ML, Deshpande A, Eaton JW, et al. Mapping male
 circumcision for HIV prevention efforts in sub-Saharan Africa. BMC Med. 2020;18:189.

17. Freeman EE, Weiss HA, Glynn JR, Cross PL, Whitworth JA, Hayes RJ. Herpes simplex virus 2 infection
increases HIV acquisition in men and women: systematic review and meta-analysis of longitudinal
studies. AIDS. 2006;20:73–83.

- 18. Chen L, Jha P, Stirling B, Sgaier SK, Daid T, Kaul R, et al. Sexual Risk Factors for HIV Infection in Early
 and Advanced HIV Epidemics in Sub-Saharan Africa: Systematic Overview of 68 Epidemiological Studies.
 PLoS ONE. 2007;2:e1001.
- 845 19. Glynn JR, Caraël M, Auvert B, Kahindo M, Chege J, Musonda R, et al. Why do young women have a
 846 much higher prevalence of HIV than young men? A study in Kisumu, Kenya and Ndola, Zambia. AIDS.
 847 2001;15:S51.
- 20. Johnson K, Way A. Risk Factors for HIV Infection in a National Adult Population: Evidence From the
 2003 Kenya Demographic and Health Survey. JAIDS J Acquir Immune Defic Syndr. 2006;42:627–36.
- 21. Pettifor AE, van der Straten A, Dunbar MS, Shiboski SC, Padian NS. Early age of first sex: a risk factor
 for HIV infection among women in Zimbabwe. AIDS Lond Engl. 2004;18:1435–42.
- 22. Coffee MP, Garnett GP, Mlilo M, Voeten HACM, Chandiwana S, Gregson S. Patterns of Movement
 and Risk of HIV Infection in Rural Zimbabwe. J Infect Dis. 2005;191 Supplement_1:S159–67.
- 854 23. Bwayo J, Plummer F, Omari M, Mutere A, Moses S, Ndinya-Achola J, et al. Human Immunodeficiency
 855 Virus Infection in Long-Distance Truck Drivers in East Africa. Arch Intern Med. 1994;154:1391–6.

24. Davis KR, Weller SC. The effectiveness of condoms in reducing heterosexual transmission of HIV. Fam
Plann Perspect. 1999;31:272–9.

- 25. Duncan ME, Peutherer JF, Simmonds P, Young H, Tibaux G, Pelzer A, et al. First coitus before
 menarche and risk of sexually transmitted disease. The Lancet. 1990;335:338–40.
- 26. Stöckl H, Kalra N, Jacobi J, Watts C. Is Early Sexual Debut a Risk Factor for HIV Infection Among
 Women in Sub-Saharan Africa? A Systematic Review. Am J Reprod Immunol. 2013;69:27–40.
- 27. Carswell JW, Lloyd G, Howells J. Prevalence of HIV-1 in east African lorry drivers. AIDS. 1989;3:759–
 62.
- 28. Brown NR, Sinclair RC. Estimating number of lifetime sexual partners: Men and women do it
 differently. J Sex Res. 1999;36:292–7.
- 29. Golding N, Burstein R, Longbottom J, Browne AJ, Fullman N, Osgood-Zimmerman A, et al. Mapping
 under-5 and neonatal mortality in Africa, 2000–15: a baseline analysis for the Sustainable Development
 Goals. The Lancet. 2017;390:2171–82.
- 30. Rue H, Martino S, Chopin N. Approximate Bayesian inference for latent Gaussian models by using
 integrated nested Laplace approximations. J R Stat Soc Ser B Stat Methodol. 2009;71:319–92.
- 31. GADM database of global administrative areas. 2018. https://gadm.org/ (accessed May 6, 2018).
 Accessed 6 May 2018.
- 873 32. Breiman L. Stacked regressions. Mach Learn. 1996;24:49–64.
- 874 33. Bhatt S, Cameron E, Flaxman SR, Weiss DJ, Smith DL, Gething PW. Improved prediction accuracy for
 875 disease risk mapping using Gaussian process stacked generalization. J R Soc Interface. 2017;14.
- 876 34. Wood SN. Generalized additive models: an introduction with R, second edition. CRC Press; 2017.
- 877 35. Chen T, He T. xgboost: eXtreme gradient boosting. :4.
- 36. Friedman JH, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate
 descent. J Stat Softw. 2010;33:1–22.
- 37. Kuhn M, Wing J, Weston S, Williams A, Keefer C, Engelhardt A, et al. caret: classification and
 regression training. 2020.
- 38. Marquez N, Wakefield J. Harmonizing child mortality data at disparate geographic Levels.
 ArXiv200200089 Stat. 2020.
- 39. Gouws E, Mishra V, Fowler TB. Comparison of adult HIV prevalence from national population-based
 surveys and antenatal clinic surveillance in countries with generalised epidemics: implications for
 calibrating surveillance data. Sex Transm Infect. 2008;84 Suppl 1:i17–23.
- 40. Marsh K, Mahy M, Salomon JA, Hogan DR. Assessing and adjusting for differences between HIV
 prevalence estimates derived from national population-based surveys and antenatal care surveillance,
 with applications for Spectrum 2013. AIDS Lond Engl. 2014;28:S497–505.

- 41. Kongnyuy EJ, Wiysonge CS. Association between fertility and HIV status: what implications for HIVestimates? BMC Public Health. 2008;8:309.
- 35. Zaba, B. & Gregson, S. Measuring the impact of HIV on fertility in Africa. AIDS12 Suppl 1, S41-50 (1998).
- 43. GBD 2019 Diseases and Injuries Collaborators. Global burden of 369 diseases and injuries in 204
 countries and territories, 1990–2019: a systematic analysis for the Global Burden of Disease Study 2019.
 Lancet. 2020;396:1204–22.
- 44. Giorgi E, Sesay SSS, Terlouw DJ, Diggle PJ. Combining data from multiple spatially referenced
 prevalence surveys using generalized linear geostatistical models. J R Stat Soc Ser A Stat Soc.
 2015;178:445–64.
- 45. Diggle PJ, Giorgi E. Model-Based Geostatistics for Prevalence Mapping in Low-Resource Settings. J
 Am Stat Assoc. 2016;111:1096–120.
- 46. Franco-Villoria M, Ventrucci M, Rue H. A unified view on Bayesian varying coefficient models.
 ArXiv180602084 Stat. 2019.
- 47. Simpson D, Rue H, Riebler A, Martins TG, Sørbye SH. Penalising Model Component Complexity: A
 Principled, Practical Approach to Constructing Priors. Stat Sci. 2017;32:1–28.
- 48. Fuglstad G-A, Simpson D, Lindgren F, Rue H. Constructing Priors that Penalize the Complexity of
 Gaussian Random Fields. J Am Stat Assoc. 2019;114:445–52.
- 49. Kristensen K, Nielsen A, Berg CW, Skaug H, Bell BM. TMB: Automatic Differentiation and Laplace
 Approximation. J Stat Softw. 2016;70:1–21.

50. Lindgren F, Rue H, Lindström J. An explicit link between Gaussian fields and Gaussian Markov
random fields: the stochastic partial differential equation approach. J R Stat Soc Ser B Stat Methodol.
2011;73:423–98.

- 51. Murray CJ, Ezzati M, Flaxman AD, Lim S, Lozano R, Michaud C, et al. GBD 2010: design, definitions,
 and metrics. The Lancet. 2012;380:2063–6.
- 52. Roberts DR, Bahn V, Ciuti S, Boyce MS, Elith J, Guillera-Arroita G, et al. Cross-validation strategies for
 data with temporal, spatial, hierarchical, or phylogenetic structure. Ecography. 2017;40:913–29.
- 53. Waller LA. Estimate suggests many infant deaths in sub-Saharan Africa attributable to air pollution.
 Nature. 2018;559:188–9.
- 919 54. Burstein R, Henry NJ, Collison ML, Marczak LB, Sligar A, Watson S, et al. Mapping 123 million
- neonatal, infant and child deaths between 2000 and 2017. Nature. 2019;574:353–8.