

METHODOLOGY FOR THE

# Unit Cost Study Repository

Lily Alexander, Lori Bollinger,  
Drew Cameron, Lauren Carroll, Willyanne DeCormier Plosky, Gabriela Gomez, Carol Levin, Mariana Siapka

**CHCC**

METHODOLOGY FOR THE

# UNIT COST STUDY REPOSITORY

Lily Alexander, Lori Bollinger, Drew Cameron,  
Lauren Carroll, Willyanne DeCormier, Gabriela Gomez,  
Carol Levin, Mariana Siapka

# CONTENTS

ACRONYMS -----Error! Bookmark not defined.

INTRODUCTION ----- 1

OBTAINING THE PUBLISHED DATA -----2

    Systematic search and review -----2

    Development of the extraction form -----9

    Data extraction and cleaning process ----- 13

CLEANING AND PROGRAMMING THE DATA ----- 15

    Generating the wide file from Excel----- 15

    Validation of data transformations----- 16

    Challenges and insights----- 19

    Importing the wide file into the UCSR----- 19

THINGS TO WATCH FOR----- 21

    What is to come ----- 21

    Limitations ----- 21

# ACRONYMS

ART	Antiretroviral Therapy
DCP3	Disease Control Priorities Project
EAGLE	European Association for Grey Literature Exploitation
GHCC	Global Health Cost Consortium
HIV	Human Immunodeficiency Virus
IEC	Information, Education, and Communication
LILACS	Literatura Latinoamericana en Ciencias de la Salud
LSHTM	London School of Hygiene and Tropical Medicine
OI	Opportunistic Infections
PLHIV	People Living with HIV
PMTCT	Prevention of Mother-to-Child Transmission
SIGLE	System for Information on Grey Literature in Europe
TB	Tuberculosis
UCSR	Unit Cost Study Repository
USD	U.S. Dollars
VMMC	Voluntary Male Medical Circumcision
WHO	World Health Organization

# INTRODUCTION

A critical gap in the arsenal needed for planning tuberculosis (TB) and HIV programs is a centralized source of standardized intervention cost data that is easily accessible to policy analysts, country officials and implementing organizations. Access to accurate intervention costs that are relevant to local contexts could help support the costing of national strategies, assist in Global Fund applications, identify opportunities for sustainability, and perform economic evaluations, including identifying potential inefficiencies. The Unit Cost Study Repository (UCSR) gathers together in one easily accessible online platform all published and grey literature cost estimates for TB and HIV interventions. The cost estimates have been standardized by the Global Health Cost Consortium (GHCC), in consultation with expert advisors, stakeholders, and partners, in terms of output units (e.g., per person served, per visit), intervention implementation (e.g., service delivery platforms, ownership, target populations, technologies), disaggregated cost categories (e.g., personnel, capital costs), and costing perspectives.

The UCSR is designed to be easy to use, whether on a desktop or on a mobile device. Costs are primarily categorized by intervention, and users may choose to display available cost data for specific interventions after filtering for the disease (HIV, TB) and intervention class (prevention, treatment, etc.). The list of interventions, and the categorization of the interventions by intervention class, align with a standardized typology of interventions, developed after extensive consultation with partners. It defines the scope of TB and HIV interventions and classifies them in a manner consistent with the monitoring and budgeting structures of GHCC partners. Thus, a user searching for cost data to fill in a Global Fund application, or to evaluate (from the donor side) the costs provided in an application, can use the UCSR to quickly find cost data that aligns with the intervention definitions and classifications in their own reporting structures.

The UCSR has additional filters for geography (region, country, urbanicity), target population (demographic, clinical), and implementation (platform, ownership, technology) so that the user may narrow the search to find locally relevant cost data. In addition, the user may also refine the search by costing methodology (perspective, economic/financial cost, etc.), and may view detailed information about how the intervention and study were conducted in secondary pages (e.g., staff type and numbers, year of the cost data collection, discount rate). This allows the user to clarify why one cost estimate may be higher or lower than other estimates that are displayed. Every attempt was made to provide clarity to users regarding cost differences across studies for the same intervention, while maintaining ease of use.

# OBTAINING THE PUBLISHED DATA

## Systematic search and review

In order to populate the UCSR, the GHCC performed systematic reviews to identify empirical cost articles and reports for HIV and TB interventions among low- and middle-income countries. The systematic review process is first described below for HIV, followed by a description of the process followed for TB.

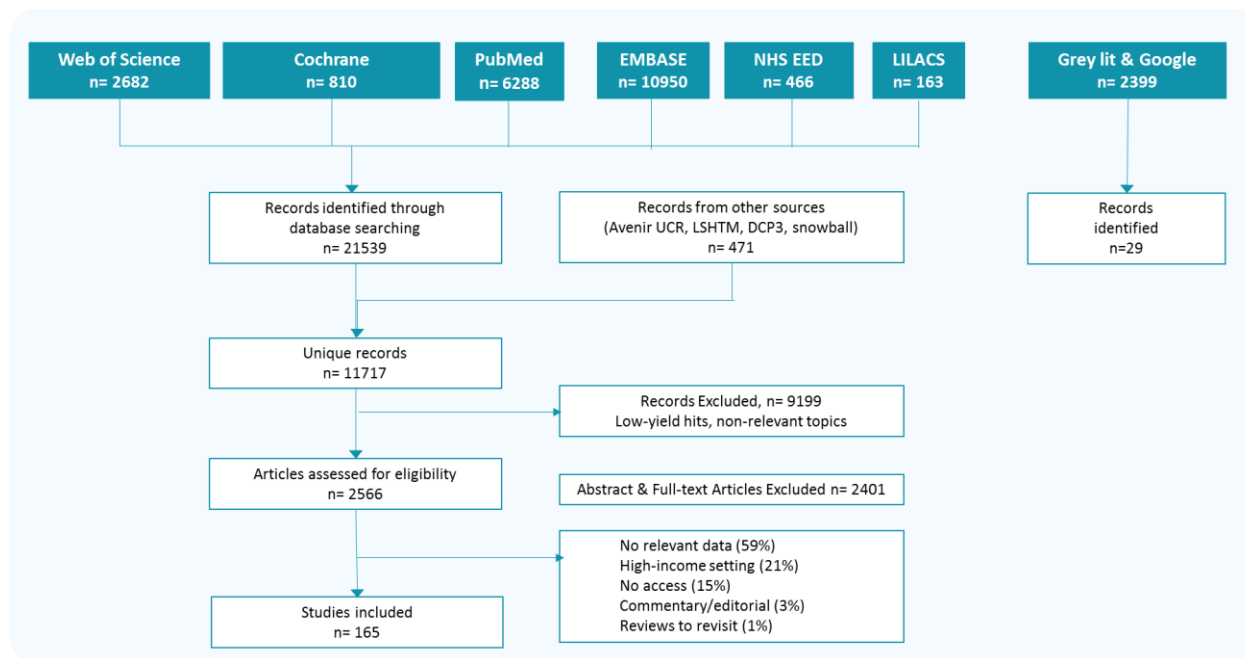
**HIV systematic search and review:** GHCC collaborated with partners to identify 33 HIV interventions covering prevention, treatment and care, testing, enablers, and health systems. We conducted a systematic review of published literature between January 2006 and October 2016, using a combination of search terms tailored to each specific database which generally combined cost terms with HIV disease terms (see Table 1). There were no restrictions on the types of treatment or interventions for HIV and AIDS, as the goal was to capture as many relevant studies as possible. The initial search was limited to articles published from January 1, 2006 through the search date October 2, 2016. Further, a search with the same databases was completed on October 20, 2017 (date range from January 1, 2006 to October 20, 2017) to identify emerging studies on differentiated care models for ART delivery. To supplement the above searches, we used databases from systematic reviews previously conducted by Avenir Health (Unit Cost Repository), London School of Tropical Hygiene and Medicine (LSHTM), and the Disease Control Priorities Project (DCP3). In addition, we conducted a search of the grey literature and focused Google searches to capture as many relevant studies as possible.

TABLE 1 – HIV SEARCH PARAMETERS AND DATABASES

Databases	Search terms
PubMed	<p>1. HIV Infections[MeSH] OR HIV[MeSH] OR hiv[tw] OR hiv-1*[tw] OR hiv-2*[tw] OR hiv1[tw] OR hiv2[tw] OR hiv infect*[tw] OR human immunodeficiency virus[tw] OR human immunodeficiency virus[tw] OR human immuno-deficiency virus[tw] OR human immune-deficiency virus[tw] OR ((human immun*) AND (deficiency virus[tw])) OR acquired immunodeficiency syndrome[tw] OR acquired immunodeficiency syndrome[tw] OR acquired immuno-deficiency syndrome[tw] OR acquired immune-deficiency syndrome[tw] OR ((acquired immun*) AND (deficiency syndrome[tw])) OR "Sexually Transmitted Diseases, Viral"[MeSH:noexp]</p> <p>2. cost*[Title/Abstract] OR "costs and cost analysis"[MeSH:noexp] OR cost benefit analys*[Title/Abstract] OR cost-benefit analysis[MeSH Term] OR health care costs[MeSH:noexp]</p> <p>3. #1 AND #2</p>

Embase	<ol style="list-style-type: none"> <li>1. 'hiv'/mj OR 'aids'/mj OR 'human immunodeficiency virus'/exp OR 'acquired immunodeficiency syndrome'/exp OR 'hiv':ab,ti</li> <li>2. 'health care cost'/exp/mj OR 'health care cost' OR 'cost'/exp OR cost OR costs</li> <li>3. #1 AND #2</li> </ol>
Web of Science	<ol style="list-style-type: none"> <li>1. "Health care cost" OR "health care costs" or cost or costs</li> <li>2. "Human immunodeficiency virus" OR "Human immune deficiency virus" OR "acquired immunodeficiency syndrome" OR "acquired immune deficiency syndrome" OR HIV* OR HIV/AIDS OR AIDS</li> <li>3. #1 AND #2</li> </ol>
Cochrane Central Register of Controlled Trials; Cochrane Reviews	<ol style="list-style-type: none"> <li>1. "Health care cost" OR "health care costs" or cost or costs</li> <li>2. "Human immunodeficiency virus" OR "Human immune deficiency virus" OR "acquired immunodeficiency syndrome" OR "acquired immune deficiency syndrome" OR HIV* OR HIV/AIDS OR AIDS</li> <li>3. #1 AND #2</li> </ol>
NHS Economic Evaluations Database, via Cochrane Library	<ol style="list-style-type: none"> <li>1. "Human immunodeficiency virus" OR "Human immune deficiency virus" OR "acquired immunodeficiency syndrome" OR "acquired immune deficiency syndrome" OR HIV* OR HIV/AIDS OR AIDS</li> </ol>
Literatura Latinoamericana en Ciencias de la Salud (LILACS)	<ol style="list-style-type: none"> <li>1. HIV OR VIH OR HIV/AIDS OR VIH/SIDA OR AIDS OR SIDA</li> <li>2. cost OR costs OR costo OR costos OR custo OR custos</li> <li>3. #1 AND #2</li> </ol>

Results from the systematic search were stored in an EndNote Library and merged with article lists previously obtained by LSHTM, DCP3, and Avenir Health, along with literature obtained from snowball searches. For the published literature, the 21,539 records identified in the systematic search were merged with 471 records from these other sources, resulting in 11,717 unique records (see Figure 1). An initial screen excluded studies with irrelevant topics (e.g., animal or in-vitro, hearing aids, high-income country settings), bringing the total down to 2,566 studies. A team of four researchers screened the results and, based on title and abstract, articles were excluded for the following reasons: high-income settings, no empirical cost data (e.g., modeled studies, program evaluation studies, etc.), commentary or editorial article, lack of access, or review articles. Senior researchers completed random checks of the excluded studies to identify potentially missed studies. All articles and reports were screened a second time during the extraction process to ensure that they contained the desired empirical cost data. A final set of 165 peer-reviewed articles was extracted.

**Figure 1: PRISMA Diagram for HIV**

A search of the grey literature, including focused Google searches, yielded a total of 2,399 potential reports. The grey literature search focused on relevant websites, while the Google searches focused on identifying specific interventions by region. After excluding duplicates and studies not meeting inclusion criteria, the grey literature search resulted in 25 reports. Focused Google searches were completed between November 2017 and March 2018 with the following search string format: “[intervention name]” costs [one of: Africa, Asia, East Europe]-US. For example, a search string for female condom provision would be: “female condom provision” costs Africa-US. The first 20 results for each region were reviewed by title and abstract, resulting in another four unique reports for inclusion.

Note that an updated systematic search is underway to identify recent HIV-related publications for inclusion in the UCSR. Table 2 displays the number of studies by intervention, as well as the corresponding number of unit cost estimates by intervention. Note that the latter number may be higher, as many studies contain more than one unit cost estimate.

**TABLE 2 – NUMBER OF STUDIES AND UNIT COSTS BY INTERVENTION FOR HIV**

Intervention	Article/reports	Unit cost estimates
Information, Education, and Communication (IEC)	6	32
HIV Counseling and Testing	26	71
Male Condom Provision	1	1
Female Condom Provision	1	1



Cash Transfers	0	0
Service Package for Key Population	15	42
Needle and Syringe Programs	5	14
Opioid Substitution Therapy	6	12
STI Management	8	38
Blood Safety	1	3
PEP	1	2
Injection Safety	0	0
VMMC	32	93
Pre-Exposure Prophylaxis (PrEP)	1	5
PMTCT	9	61
Adult ART	61	208
Differentiated Care ART	1	4
Pediatric ART	12	50
Linkage to Care	1	2
Retention and Adherence	2	37
Inpatient Care	7	37
Post-violence Care	1	4
Patient Tracking	1	2
Condom Social Marketing (CSM)	1	2
Workplace Service Package	4	76
OI Prophylaxis	2	2
OI Diagnosis and Treatment	2	4
Socioeconomic Support for PLHIV	4	85
Laboratory Monitoring	3	13
Pre-ART Care	8	15
HIV/TB Care Delivery	3	38
Provider Engagement/Training	1	1
Supply Chain Management	3	29
Infection Control	0	0
Stigma Reduction and Human Rights	0	0
Health System Interventions	0	0
Stigma Reduction	0	0
Community Empowerment	0	0
Gender-Based Violence Prevention	0	0
Income Generation	0	0
In-kind Benefits (e.g., clothing)	0	0
Counseling and Psychosocial Support for PLHIV	0	0
Legal Literacy	0	0
Training Providers on Human Rights	0	0
Legal Services	0	0

Human Rights Legislation	0	0
Monitoring and Evaluation Systems	0	0
Monetary Incentives for Human Resources	0	0
Training and Education	0	0
Workforce Retention and Scale-up	0	0
Planning, Coordination, and Program Management	0	0
Drug-resistance Surveillance	0	0
HIV Serological Surveillance	0	0
Grant Management and Disbursement	0	0
Operations Research/Quality Improvement	0	0

**TB systematic search and review.** For TB, six health and economic databases were searched between May and July 2016: Pubmed, EMBASE, Econlit, The National Health Service Economic Evaluation Database, The Cost-effectiveness Analysis Registry, and Cochrane library. Two additional databases – Web of Science and Latin American and Caribbean Health Sciences Literature (LILACS) – were searched in February and March 2017.

Broad searches were designed with search terms including disease, intervention, and outcome of interest only (see Table 3). We searched for all TB interventions, ensuring search terms incorporating, among other terms, drug resistant TB, prevention, and health system factors.

TABLE 3 – TB SEARCH TERMS

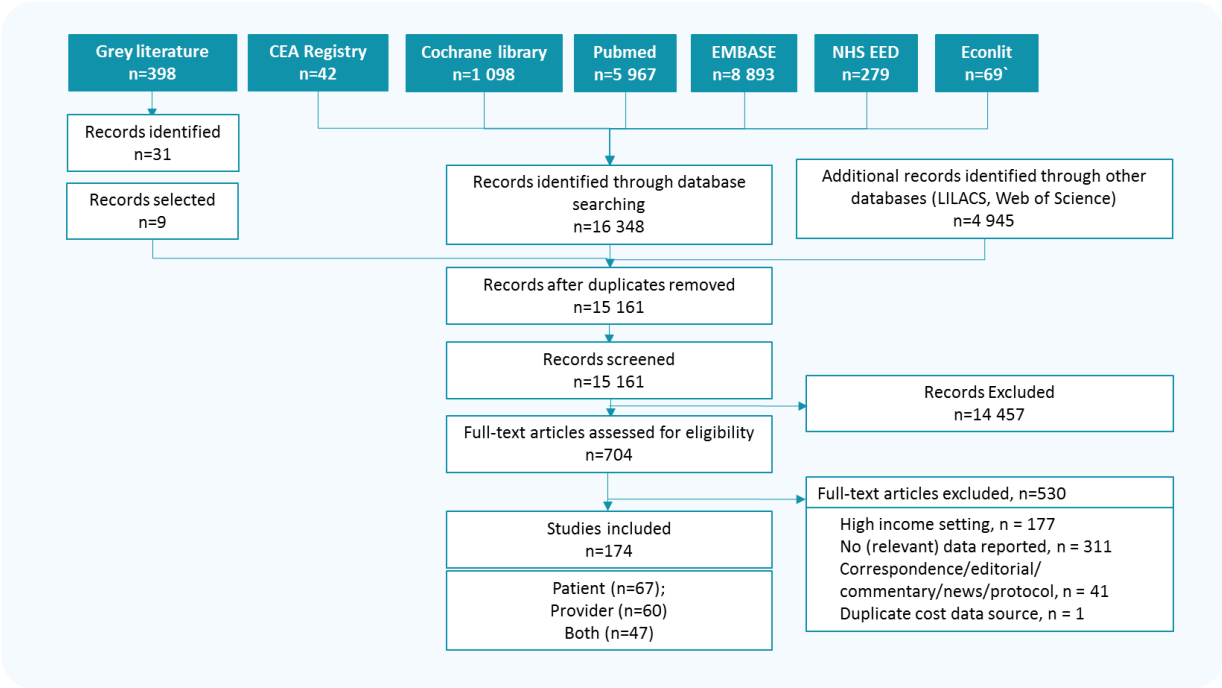
Category	Search terms
1. Cost	Cost* or economic or finance AND
2. Tuberculosis	TB or tuberculosis or MDR#TB or XDR#TB or multi?drug or “resistant tuberculosis” or “strain resistance” or “mycobacterium tuberculosis” AND
3. Treatment	treatment or management or drugs or medication or DOTS or “directly observed treatment” or “health system*” or “hospital care” or “epidemiology” or “government hospital setting” or “community based care” or “patient* perspective” or “isoniazid preventive therapy” or “IPT” or “prevention”

We included all studies reporting TB unit or output costs published between January 1990 and July 28, 2016, without any restriction on language or geography. We excluded articles if they had no empirically collected data or if we could not ascertain the currency or perspective of the costing reported (from the article or after contacting authors).

The search results from each database were downloaded to Endnote. Overall, a total of 15,161 articles were identified after excluding duplicates (see Figure 2 for results by database).

The screening of these results was then undertaken in three stages: by title, by abstract, and by full text. First, one person reviewed all records by title, while a second person checked the records that were excluded based on title screening. Records were excluded at this stage if it was evident from the title that the study referred to research in animals. Second, the resulting 6,307 abstracts were screened by two people independently, and were excluded if they did not report costs or costing. This resulted in a final number of 704 studies assessed for eligibility for reviewing the full text ('best bets').

Figure 2: PRISMA Chart for TB



The full texts of the 704 best bets were screened by two people independently and assessed for eligibility. Records were excluded if the studies were conducted in high-income countries (n=177), no relevant data were reported (e.g., no empirical data collected regarding prices or quantities or TB-related costs reported, n=311), correspondence/editorials/commentaries/news/protocols (n=41), or contained duplicate cost data (n=1). We also cross-checked records against a recent systematic literature review<sup>1</sup> of tuberculosis costs for health services and patients.

In addition to the published literature, we searched the grey literature. We aligned our grey literature searches with the HIV team described above and searched the following sources for relevant articles, using the same exclusion criteria as was used for the peer-reviewed literature:

- The European Association for Grey Literature Exploitation (EAGLE)

<sup>1</sup> Laurence YV, Griffiths UK, Vassall A. "Costs to Health Services and the Patient of Treating Tuberculosis: A Systematic Literature Review." *Pharmacoeconomics*. 2015 Sep;33(9):939-55. doi: 10.1007/s40273-015-0279-6.

- The System for Information on Grey Literature in Europe (SIGLE)
- Documents and meeting reports from the World Bank and WHO websites

In addition, we conducted Google searches, and reviewed the first 50 documents that resulted from the algorithm used in Google for different websites (see Table 4):

TABLE 4 – TB GREY LITERATURE SEARCHES				
Web page	Organization	Search algorithm in Google	Documents reviewed	Notes
msf.org	Medecins Sans Frontieres	site: msf.org ((TB OR Tuberculosis) AND unit cost) filetype:pdf	50	Also searched the website using the terms "tb" and "cost"
who.int	World Health Organization	site: who.int ((TB OR Tuberculosis) AND "unit cost") filetype:pdf	50	Also searched the website using the terms "tb" and "cost"
unaids.org	UNAIDS	site: unaids.org ((TB OR Tuberculosis) AND cost) filetype:pdf	50	Also searched the website using the terms "tb" and "cost"
pepfar.gov	President's Emergency Plan for AIDS Relief	site: pefar.gov ((TB OR Tuberculosis) AND unit cost) filetype:pdf	50	Also searched the website using the terms "tb" and "cost"
cdc.gov	Center for Disease Control and Prevention	site: cdc.gov ((TB OR Tuberculosis) AND unit cost) filetype:pdf	50	Also searched the website using the terms "tb" and "cost"
avenirhealth.org	Avenir Health	site: avenirhealth.org ((TB OR Tuberculosis) AND cost) filetype:pdf	32 (no more available)	Also searched the website using the terms "tb" and "cost"
usaid.gov	United States Agency for International Development	site: usaid.gov ((TB OR Tuberculosis) AND unit cost) filetype:pdf	50	Also searched the website using the terms "tb" and "cost"
healthpolicyproject.com	Health Policy Project	site: healthpolicyproject.com ((TB OR Tuberculosis) AND cost) filetype:pdf	50	Also searched the website using the terms "tb" and "cost"
pangaeaglobal.org	Pangaea Global AIDS	site: pangaeaglobal.org ((TB OR Tuberculosis) AND cost) filetype:pdf	16 (no more available)	Also searched the website using the terms "tb" and "cost"

A total of 398 documents were screened as a result of the grey literature and Google searches, of which 31 (8%) were identified as potential documents containing TB unit cost data. Nine of these

reports (29%) were focused on TB and the remaining reported unit costs on interventions across multiple diseases. In addition, we identified the following costing tools, and checked for any additional unit cost information (see Table 5). While these inform the Standards and Methods work, they were not included in the UCSR as they did not contain unique cost information.

TABLE 5 – COSTING TOOLS REVIEWED

	Costing tools	Tool type	Publication year	Source
1	One Health tool	Tool	2013	Avenir Health
2	Global Price Reporting Mechanism	Database	2011	WHO
3	Unit Cost Repository	Database	2010	Avenir Health
4	Unit Cost Estimations WHO-CHOICE	Tool	2011	WHO

The same screening process was followed as for the peer-reviewed literature. When an abstract was not available, the executive summary was assessed. In total, 174 studies were included from the grey literature; of these, 107 studies reported costs from a provider perspective. Table 6 shows the number of studies and the number of unit costs by type of intervention, respectively. Several unit costs across multiple interventions were reported across studies.

TABLE 6 – NUMBER OF STUDIES AND UNIT COST ESTIMATES BY INTERVENTION FOR TB

	Article/reports	Unit cost estimates
Above Service Costs	2	4
Active Case Finding (ACF)	10	9
Intensified Case Finding (ICF)	8	11
Passive Case Finding (PCF)	27	53
TB Prevention	8	15
TB Treatment	64	366

Note that, because there have been very few costing studies in TB with large sample sizes, most unit costs were expected to come from studies with noneconomic primary outcomes. Thus, the searches were designed to be broad; however, this decreased the specificity of the searches and increased the number of studies identified not relevant for the UCSR. It was therefore important for the TB team to double check all excluded articles, due to their large number, to ensure they were excluded correctly.

## Development of the extraction form

To develop the early version of the extraction form, an extraction working group was created within the GHCC core team. This subgroup looked both at the extraction form that underpinned the former [Avenir Health] Unit Cost Repository, and the extensive Principles and Methods Reporting Checklist

from the GHCC Reference Case for Estimating the Costs of Global Health Services and Interventions, to see which potential fields were necessary. The extraction form was designed to serve multiple purposes: (a) key fields to be displayed in the UCSR, (b) additional fields to support a forthcoming quality index, and (c) other fields to support analysis. It was recognized at the outset that standardization was key to: (a) present a reasonable number of dropdown options in the UCSR, (b) compare data accurately, and (c) avoid errors that would impact programming (e.g., using Republic of Tanzania, Tanzania, Tazania). Achieving consensus across the GHCC core team, advisors, stakeholders, and partners for standardized options was a time-consuming process, and continued even after extraction began. In addition, sometimes standardized options could not be identified until several studies for a specific intervention were extracted. Extraction proceeded using a weekly data and analytics check-in call, continued discussion over email, and supplemental calls to resolve particularly tricky issues; all decisions were carefully tracked. See Figures 3 and 4 below for the final list of broad and narrow input cost categories and activity cost categories, along with descriptions:

**Figure 3: Broad and narrow Input cost categories**

BROAD AND NARROW INPUT COST CATEGORIES		DESCRIPTION OF INPUTS IN INPUT COST CATEGORIES
Personnel	Service delivery personnel	Doctors, nurses, counselors; Pharmacists; Lab/diagnostic personnel; Outreach workers, peer supporters, social workers; Community volunteers, or home visitors
	Support personnel	Administrators, supervisors; Procurement officers, supply clerks, accountants; Legal staff; Receptionists; Social media coordinators, community strategy/mobilization supervisors; Data and IT staff; Drivers; Gardeners; Security guards; Kitchen staff; Custodians or cleaning staff.
Capital	Lab/ diagnostic equipment	Centrifuges, incubators, microscopes, water baths, orbital shakers; Xpert, X-ray, microscopy instruments, GeneChip scanner; hemoglobin meters, urine analyzers, liver/renal biochemistry analyzers.
	Equipment (medical/intervention, excl. lab)	Refrigerators, freezers; Incinerators and autoclaves; MEMS caps, monitoring equipment; Tents.
	Equipment (non-medical/intervention)	Furniture: beds, benches/couches, chairs, desks, tables, lamps/fixtures, filing/drug cabinets, bookcases; Computers, monitors, LCD projectors, printers; Software; Power outlets, or paper shredders.
	Vehicles, capital	Bicycles; Motorcycles; Cars, vans or SUVs; Trucks; Boats; or Airplanes.
	Building/Space, capital	Construction/purchased floor space in a health facility or training school; Truck containers; Storage facilities; Administrative offices; Wells; or Latrines.
	Other capital	<i>Start-up training and materials; Licenses/copyrights.</i>
Recurrent	Supplies (key drugs)	TB drugs; PrEP; ARVs; PEP; Hepatitis/STI/OI Medication; Antibiotics; or Contraceptives.
	Supplies (medical/intervention, excl. key drugs)	Vaccines; Syringes, test kits, sputum bottles, speculum, cotton swabs, microscope slides reagents; Gloves, gowns, masks, bandages; Small medical equipment; or Small containers to hold drugs.
	Supplies (non-medical/non-intervention)	Pens, pencils, dry-erase markers, highlighters; printer paper, post-it notes, notebooks, calendars; paper clips, binder clips; file folders;

		envelopes, stamps; tape, glue; scissors, staplers, hole-punchers, calculators; memory sticks; batteries; or Lanyards.
	Building/space	<u>Rent</u> for capital inputs; Maintenance: Painting, roof, heating/plumbing, windows; Tires, spare parts, oil/lubricants, tune-ups; or Computer repair. Lighting, heating, water; Telephone, or internet.
Recurrent	Other recurrent	Gasoline, fuel; Tolls; or Contracted transportation services; Food (at facilities/meetings; for nutritional support to improve health or lessen side effects of drugs); Vitamins, or Contracted meal services. <i>Recurrent training</i> ; Medical malpractice insurance; Insurance for capital building, vehicles, or equipment; Registration fees for capital items, for memberships in professional organizations, or for use of copyrighted materials for communication purposes (icons, photos, etc.); Contracted services such as laboratory, storage, waste removal (even if just burning and/or burying), security, or information technology if outsourced; Courier/UPS service; or Other recurrent costs.

**Figure 4: Broad and narrow Activity cost categories**

BROAD AND NARROW ACTIVITY COST CATEGORIES		DESCRIPTION OF EXAMPLE INPUTS IN ACTIVITY CATEGORIES
Primary service delivery	Key activity 1: e.g., Voluntary medical male circumcision procedure	Doctor, nurse; Disposable surgical kit, gloves, mask, gown
	Key activity 2: e.g., Post-procedure check-up	Nurse; Gloves; Antibiotic cream
Secondary service delivery	Secondary activity 1: e.g., HIV testing and counseling	Nurse; Antiseptic, cotton pad, needle, collection tube, HIV rapid test, bandages
	Secondary activity 2: e.g., Provision of condoms	Nurse; condoms
Ancillary service delivery	Demand generation	Communication coordinator, tech/web designer; Facebook ads, radio airtime
	Lab services	Lab service fee
	Adherence/retention	Cost per text message sent to remind clients of appointments, and to remind clients to use condoms during the healing period
Operational	Buildings and equipment	
	Logistics	
	Supervision	
	Training	
	Transportation	
	Mass education	
	HMIS and record keeping	
	Technology development	

Technology maintenance	
Project management	

Following the pilot extraction of Voluntary Male Medical Circumcision (VMMC) and ART interventions, the extraction team reviewed all the fields and definitions in the extraction form to identify areas in need of standardization or editing. The changes primarily fell into three categories: fields edited to have standardized content; fields no longer included in the UCSR but retained in the extraction form for use when downloaded; and fields dropped from the extraction form.

Several fields were standardized to better support UCSR display and analysis. For example, the author-reported ‘Omitted Costs’ field was standardized to best capture any omitted costs through the use of a checklist (derived from ‘Standardized Inputs’). Further, author-reported country names were standardized into ISO-3 codes to ensure correct spelling. Additionally, the open-text field ‘Population’ was disaggregated and standardized into ‘Target Group – Demographic’ and ‘Target Group – Clinical’ to more clearly note different population groups across studies. The list of options for these fields were identified in a post-hoc analysis of author-reported information. Finally, the ‘Intervention Details’ and ‘Technology’ fields were disaggregated from a free-text field to separate intervention characteristics to allow for future searching capabilities in the UCSR.

Because many of the initial methodology fields captured in the extraction form contained detail and subtlety beyond the scope of the UCSR, they were either adapted or dropped from the extraction form. For example, ‘Lead Author’ replaced ‘Reference Author’, and ‘Number of Direct Observations’ and ‘Unit of Observation’ were dropped in favor of ‘Number of Sites’. After lengthy discussion with the team, the ‘Integrated Services’ field was dropped altogether due to challenges with extracting meaningful information in this format, and the ‘Full/Incremental Cost’ field was dropped due to challenges finding a consistent definition. Lastly, the field ‘traded vs. non-traded’ was dropped because trade status is country-dependent and thus the field is often noted as “mixed,” which has little analytical value. This process was particularly helpful in streamlining the extractions and providing more consistency in reporting across interventions.

The extraction form and process have several advantages and disadvantages. The extraction template is comprehensive and flexible, while offering opportunities for standardizing author-reported content. However, the extraction team faced a variety of challenges throughout the extraction process. With multiple extractors, consistency across extractors required constant communication and frequent revisions of the template before reaching a ‘stable’ form that minimized possible extraction differences. Moreover, the process of standardizing fields often required ad-hoc review of author-reported content by intervention, necessitating ongoing review of each intervention for themes and commonalities. As a result, version control became a priority to update older interventions that were extracted prior to a given change. Over time, this process also became more standardized, and easier to manage with the ‘stable’ extraction form. However, understanding which fields needed additional clarity would have been impossible to determine without going through the

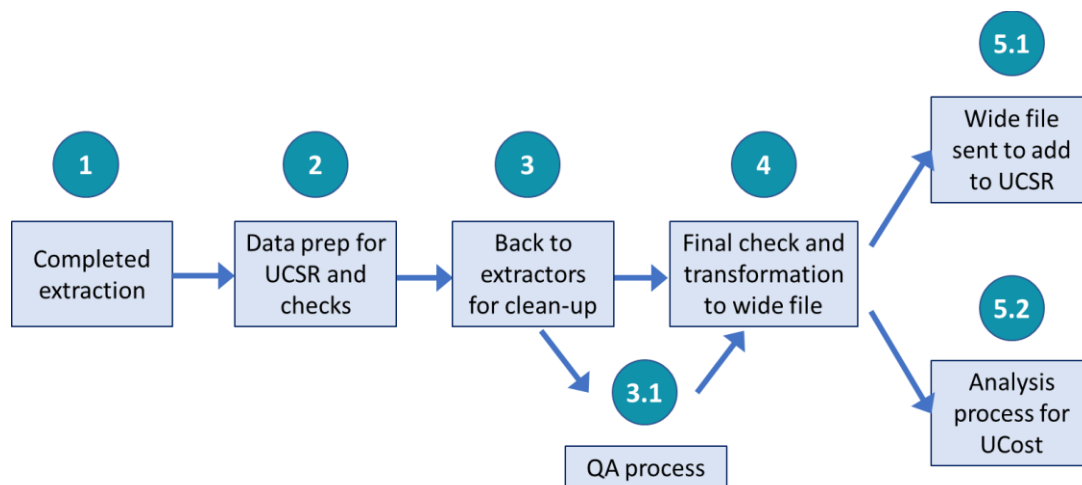


extractions. This also underscores the importance of data checks and quality-assurance steps, further described in the Data Cleaning Process section.

## Data extraction and cleaning process

The HIV data set of literature currently includes 229 articles and reports across 32 HIV interventions. Articles and reports were divided among up to three extractors, and, after initial training, each article took approximately three hours to extract. The lead extractor also merged extractions completed by different extractors for one intervention, and managed version control across all interventions. The TB data set currently includes 119 articles and reports across 6 TB interventions. Double extraction was performed by two teams of two extractors each, with significant interaction among the extractors. The following steps highlight the workflow from the initial extraction to the final upload to the UCSR (see Figure 5):

1. Once extraction for a specific intervention is completed, the extractions are reviewed for formatting issues, missing values, etc. A summary report for each intervention highlights fields needing review by the team for standardization (e.g., population and technology fields, preferred output unit costs, etc.). These summaries are reviewed for consensus on standardized options, and then the lead extractor updates the corresponding fields in the extraction.
2. The vetted interventions are sent out for data validation checks in STATA (described further below). These checks flag inconsistencies, misalignment between fields that should be aligned, and errors in cost totals. A report is sent back to the lead extractor with flags for review. These checks typically take about a day to run.
3. The extractors split up the extraction flags based on their respective familiarity with a given intervention. Depending on the number of flags, the clean-up process takes about a day. The flags are cleaned up and documented in a separate file. The cleaned-up extraction form is sent back for a final check prior to processing for the UCSR.
  - 3.1 In a parallel process, the cleaned interventions from #3 undergo quality assurance (QA) by a senior researcher. This is an iterative process that runs parallel to the other steps in the process. When completed, the data in UCSR and UCost are updated. QA takes approximately 1.25 hours per article.
4. The data validation checks are re-run, and any remaining flags are addressed by the extractors. Once cleaning is complete, the extraction file is transformed into the wide file used for UCSR and UCost.
5. The wide file is sent to Avenir Health for a final overview and incorporation into the UCSR. The file is also sent to our analytics team for use in analyses supporting UCost.

**Figure 5: Workflow for extraction to UCSR/UCost**

# CLEANING AND PROGRAMMING THE DATA

## Generating the wide file from Excel

After the completion of data extraction from published costing studies into the extraction template, data were combined into a single “wide file” for subsequent data analysis and exported to the UCSR. Before the wide file can be used for these purposes, however, cost data must be standardized to current United States Dollars (USD), and transformed data must be checked for errors in both the transformation and transcription processes. These steps are described below.

Extracted cost data were recorded on two tabs of the extraction template Excel workbook – a *Study Attributes* tab, and a *Cost Data* tab. The *Study Attributes* tab records all study-level variables, where each row represents a unique unit cost reported within a given study (e.g., urbanicity, facility type, reporting quality variables, etc.), along with corresponding characteristics in the same row. The *Cost Data* tab in the extraction template is presented in a long format for ease of transcription, recording all costs associated with each row of the *Study Attributes* tab. Each row in the *Cost Data* tab contains either an input, subtotal, or total unit cost.

Rows between the two tabs were connected by a single common identification variable formatted to include information about the disease area (e.g., ‘hiv’), the study number (e.g., ‘001’), and a counter representing the unique unit cost reported within that study (e.g., ‘a’), concatenated into a single variable during extraction (e.g., ‘hiv001a’). Further, types of interventions were grouped numerically for HIV only (e.g., ‘hiv001a’ – ‘hiv099z’ for voluntary medical male circumcision studies; ‘hiv100a’ – ‘hiv199z’ for anti-retroviral treatment interventions, etc.). In TB, studies tend to report on different types of interventions and therefore identifiers cluster around study numbers and not interventions (e.g., tb001 will have tb001a, a unit cost for diagnosis; tb001b, a unit cost for first-line treatment; and tb001c, a unit cost for MDR treatment). Thus, in the *Study Attributes* tab, each row is identified by a unique ID (e.g., hiv001a, hiv001b, hiv002a, etc.). In the cost data tab, however, each row is identified within one of these ID variables, as well as a lower case roman numeral (e.g., hiv001a\_i, hiv001b\_i, hiv002a\_i, etc.) to distinguish between other unit cost estimates such as sensitivity analyses.

In order to reshape extracted data into usable format, all data for a given unit cost reported by study authors must be included in a single row of an Excel workbook. Thus, multiple rows of input costs in the *Cost Data* tab must be combined into a single row associated with a single total unit cost, and then appended to the associated row from the *Study Attributes* tab. Figure 6 below uses hypothetical examples of extracted data to illustrate the distinction between the *Study Attributes* tab (item a) and the *Cost Data* tab (item b). Item c shows how data from these two tabs are combined in a single wide file using Stata.

## Figure 6: Examples of extraction template and wide file structures for reshaping

a. Study Attributes tab				b. Cost Data tab						
id	author	urbanicity	Country	id	Sensitivity	mean cost	broad cost	narrow cost	currency year	currency
hiv001a	Smith 2016	urban	Uganda	hiv001a	i	\$1.50	Capital	building space	2016	USD
hiv001b	Smith 2016	rural	Uganda	hiv001a	i	\$0.23	Personnel	Support staff	2016	USD
hiv002a	Bollinger 2016	mixed	South Africa	hiv001a	i	\$5.00	Personnel	Doctor salary	2016	USD
hiv002b	Bollinger 2016	rural	South Africa	hiv001a	i	\$6.73	Total	Total	2016	USD
hiv003a	Kahn 2005	urban	Tanzania	hiv001a	ii	\$7.76	Total	Total	2016	USD
				hiv001b	i	\$5.57	Personnel	doctor salary	2016	USD
				hiv001b	i	\$0.00	Personnel	Support staff	2016	USD
				hiv001b	i	\$2.51	Traded goods	circumcision kit	2016	USD
				hiv001b	i	\$0.00	Capital	Building space	2016	USD
				hiv001b	i	\$8.08	Total	Total	2016	USD
				hiv001b	ii	\$8.75	Total	Total	2016	USD
				hiv002a	i	725.15	Total	Total	2016	SA Rand
				hiv002b	i	862.55	Total	Total	2016	SA Rand
				hiv003a	i	\$27.50	Total	Total	2004	USD

c. Wide file (combined Study Attributes with Cost Data tabs in Stata)												
id	author	urbanicity	country	total	capital	personnel	traded goods	Personnel: direct service	Personnel: support	Capital: building space	Traded goods: medical	Traded goods: non-medical
hiv001a_i	Smith 2008	urban	Uganda	\$6.73	\$1.50	\$5.23	.	\$5.00	\$0.23	\$1.50	.	.
hiv001a_ii	Smith 2008	urban	Uganda	\$7.76	.	.	.	.	.	.	.	.
hiv001b_i	Smith 2008	rural	Uganda	\$8.08	\$0.00	\$5.57	\$2.51	\$5.57	\$0.00	\$0.00	\$2.51	.
hiv001b_ii	Smith 2008	rural	Uganda	\$8.75	.	.	.	.	.	.	.	.
hiv002a_i	Bollinger 2014	mixed	South Africa	\$53.87	.	.	.	.	.	.	.	.
hiv002b_i	Bollinger 2014	rural	South Africa	\$64.07	.	.	.	.	.	.	.	.
hiv003a_i	Kahn 2005	urban	Tanzania	\$35.67	.	.	.	.	.	.	.	.

Note that, before data were reshaped, costs were standardized into current USD. All costs reported in currencies other than USD were converted to USD based on the year of reporting using market exchange rates published by the World Bank. Once all costs were converted to USD, costs were inflated using the US GDP price deflator, as needed, also from the World Bank. Alternative methods of currency-year transformation were explored in great detail and are explained elsewhere.

Once costs were standardized to current year USD, data from both the *Study Attributes* and *Cost Data* tabs were imported into Stata in two separate Stata *.dta* files. Data from the *Cost Data* file were merged field-by-field into the *Study Attributes* file until each unit cost was recorded in a single row in the new wide file, with one column for the total unit cost, and subsequent columns for each broad cost category (e.g., capital, personnel, traded goods, etc.) and each narrow cost category (e.g., personnel: direct service delivery; support staff) as applicable (see Figure 6 item c for an illustrative example). The wide file contains many more categories of input costs than those listed in the figure, retaining the same structure capturing all data originally reported in the extraction template.

## Validation of data transformations

Once in a wide file format, data underwent a series of validation checks. These were done in three steps for each intervention separately using Stata 15. The overarching process for data checks and validations can be seen in Figure 5 above.

First, data were checked to ensure that costs sum correctly. This included narrow cost and activity categories summing to broad cost categories (e.g., personnel: direct service + personnel: support = personnel (broad)), as well as these broad categories summing to the total unit costs (e.g., capital +

personnel + traded goods = total unit cost). Slight deviations were expected due to cost estimates reported from authors, or modifications from exchange rate changes. All costs for which there was a discrepancy between the sum of the inputs and the reported total were flagged to be checked by the extractors.

Second, any studies that omitted key personnel, commodities, or services from their reported total unit cost were flagged for further examination. Because incomplete unit costs could be misleading to UCSR users, before transferring data to the UCSR total unit cost was treated as a missing value in cases where key inputs were missing. Decisions on when omissions warrant reporting the total as missing were determined for each intervention separately and checked during the validation process. A set of these intervention-specific determinations of omissions is documented in the Appendix.

Third, alignment between the *costed activities* and *intervention details* fields was checked. *Intervention details* fields describe how the intervention in a study was implemented (e.g., demand generation, clinical monitoring, screening and diagnoses, or treatment). The *costed activities* field identifies which elements of the intervention are included in the cost estimates. By comparing the contents of the intervention details and costed activities field, we identified discrepancies (i.e., intervention elements that were included in the intervention but not costed, or vice versa).

Variable labels and names were also checked to ensure that they were standard across interventions. Once this process was complete, flagged errors were sent back to extractors. A PDF summary of each variable in the dataset was sent as well, showing ranges for continuous variables and category labels for categorical ones (Figure 7). This is helpful for identifying anomalous values and outliers. Once extractors reviewed and addressed issues identified through the data validation process, they sent back a revised extraction form, which then underwent the reshaping and data checking process once more before it was outputted as a cleaned wide file to upload to the UCSR.

Figure 7: Summary codebook of each variable to detect anomalies and outliers

id_pop_dem_std	Target group (demographic)
type: string (str41)	
unique values: 5	missing "": 0/101
tabulation: Freq. Value	
11 "Adolescent Boys and Adult Men"	
2 "Adolescent Boys and Young Men, ages 15-29"	
5 "Infants, ages 1month - 2 yrs"	
82 "Men"	
1 "Neonates, ages birth - 1 month"	
warning: variable has embedded blanks	
id_pop_clin_std	Target group (clinical)
type: string (str12)	
unique values: 2	missing "": 0/101
tabulation: Freq. Value	
100 "HIV negative"	
1 "HiV negative"	
warning: variable has embedded blanks	

The process for TB was similar, although the steps were followed in a slightly different order. First, double extraction was performed for the extraction form, which provided for quality assurance even before exporting into the wide file. For example, one check was whether costs sum correctly; any discrepancy was checked against the publication. In addition, omitted categories were checked against the publication and standardized. Finally, unit costs missing key inputs were recorded as subtotals and have a missing value for total unit costs.

Once the Excel dataset was clean, the same two data files as mentioned above, *Study Attributes* .dta and *Cost Data* .dta were created. Using Stata 15, reporting across all fields in the two data files were checked and standardized, including whether reported units, activities and outputs matched the Reference Case, and all fields with outliers. In addition, we checked that the following categories matched: Activity category broad versus Standardized Input classification broad; Intervention category versus Intervention Type; Intervention type versus modality; Intervention type versus

technology; Intervention type versus treatment phase; and Intervention type versus costed activities. The final, merged dataset was then checked for completeness and accuracy before uploading to the UCSR.

## Challenges and insights

This process highlights several challenges in the extraction and validation of cost data. First, there is a tradeoff between listing many columns in the extraction template for input costs to be reported and doing so vertically in a different extraction template format. The vertical process is superior from an extraction standpoint because extractors can easily see and organize which costs apply to which total unit costs collected in a manner most closely mimicking the actual reporting of study authors. This limits the chances of transcription errors. Unfortunately, vertical transcription creates challenges for reshaping and combining costs into broad and narrow cost categories in a wide (or flat) file format. This in turn requires more data validation procedures to ensure that data were not lost or inaccurately summed in the reshaping process.

Second, the GHCC policy is to extract all unit costs reported by study authors (as long as they reflect real, not modeled, data). For HIV, this sometimes involved extracting duplicate data (e.g., a country average unit cost; several regional averages within a country). To distinguish between these potentially duplicate costs, the extraction team records key information about each cost in a field that is meant to help analysts and data validators to distinguish between costs. Third, there are many intervention-specific nuances to data extraction that make recording in a standard template difficult. Thus, for some interventions, intervention-specific categories within existing variables may be necessary for clarity. Going forward, one suggestion is to base extraction on the Reference Case interventions, activities, outputs, etc., in order to have standardized forms. Finally, it should be noted that having the “as reported” fields was very useful during the validation and data checking processes.

## Importing the wide file into the UCSR

After the wide file was sent to be uploaded into the UCSR, Avenir Health performed one last set of checks (e.g., existence of outliers, correctly summed variables, sensible entries, completeness of data).

Piloting the uploading of the wide file resulted in instituting a series of checks with both the extractors and those preparing the wide files in STATA. Once the system of checks was in place, Avenir Health’s data check became more focused on standardization between the TB and HIV wide files and probing where new fields should be added to support better user understanding of the cost. For example, fields on alerts (e.g., omission of key costs), number of sites included in the unit cost, and unique trait (e.g., personnel type) were added to show users detail not captured by the main characteristic fields (e.g., country, platform). This effort required continuous dialogue with the data team, and with the programmers at Avenir responsible for setting up the UCSR.

The final UCSR design incorporated input from the GHCC team, a website design firm, and a data visualization specialist to align the look of the UCSR with that of the GHCC website, and to design the “user journey” for the UCSR. The final design required the user to select one intervention (along with disease and intervention class), and allowed for further filtering by geographic location (region/country, urbanicity) and target population on the main filter page. Secondary filters appear on a “Further refine” page. After filters were selected, users chose between various data visualizations to display their results; they may also download the underlying data. The intention from the outset was not only to display a collection of cost data, but to make clear what defines each cost estimate, and to do so in a straightforward manner.



# THINGS TO WATCH FOR

## What is to come

The UCSR will be publicly launched for beta-testing at the IAEN Pre-Conference meeting in July 2018. It will also be presented at the International Union Against Tuberculosis and Lung Disease in October 2018, as well in several World Bank learning platforms. The GHCC also expects to develop a core beta-testing group with key users at the country level and within the central offices of partner organizations in late summer of 2018, who will provide valuable feedback. Finally, the UCSR will also be piloted in some country applications that align with planning cycles for Global Fund applications or national strategic plans.

The UCSR initially focuses on TB and HIV, with a view to becoming a platform integrating other health areas should the demand for the database exist. Requests to share the UCSR development and extraction process with teams working on immunization, malaria, and social behavioral change interventions have already been received over the past two years. The GHCC team has worked to provide open access to our process and templates, and we are working to develop further training materials so that teams in other health areas may adapt GHCC templates for their specific use, while maintaining the standardized options for fields, so that the data could be added to the UCSR in the future if desired.

The beta-testing will help to refine the look and functionality of the UCSR. For example, prototype data visualizations will be available in July 2018, and beta-testing will help resolve questions about the types of data and charts that users want to see. More data visualizations may be developed, and several static visualizations/infographics presenting key data from the UCSR will be developed and posted to the GHCC website. The GHCC also will add cost disaggregation in the original year of the unit cost data (as reported by authors) to the existing disaggregation displays. Finally, while the UCSR does have a translation feature utilizing the new AI Google translate, key pages (e.g., search filters, dropdown options) will be sent to official translators, with improved translations available in the fall of 2018. It is expected that the UCSR will evolve over time to meet the needs and feedback of users, particularly if other health areas are integrated.

## Limitations

Users of the UCSR are currently limited to searching by one intervention at a time. This limitation was intentional, to prevent overwhelming the user with an unmanageable number of search results, and to target the design of the data visualizations to what would be most useful to the user. Users may download the full dataset with all interventions and conduct their own external analysis and/or visualizations. For the UCSR visualizations, users are limited only to charts where it is possible to view

multiple cost estimates for the same intervention. For example, there may be six cost estimates for VMMC in Zambia, and five in Uganda; the displayed estimates would vary by either implementation features (e.g., sector), different costing methodologies, and/or different inputs included in the cost. Since the UCSR is not set up to produce an average cost for each country, only visualizations such as a box plot (showing the range of cost for each country, and the median) are possible. This is because: a) with the exception of standardization of nomenclature and inflation/currency conversion to a common year/currency, the GHCC tried to faithfully represent the study data as reported; b) the range of estimates may be overly influenced by one study that has a multitude of estimates (such as for each facility studied); and c) the UCost Tool will give an “average” point estimate for the cost of each intervention in each country, on the basis of data drawn from extensive analysis of literature and primary data. Finally, since many of the fields in the UCSR are subjective, there may be differences among extractors in how data were extracted. even with extensive training of extractors and data quality checks. Accordingly, we will be contacting the lead author of each study in the UCSR, to ask them to review and validate the information presented.