



# HHS Public Access

Author manuscript

*Am J Public Health*. Author manuscript; available in PMC 2018 July 01.

Published in final edited form as:

*Am J Public Health*. 2017 July ; 107(7): 1078–1086. doi:10.2105/AJPH.2017.303707.

## REVIEW OF RECENT METHODOLOGICAL DEVELOPMENTS IN GROUP-RANDOMIZED TRIALS: PART 2 - ANALYSIS

**Elizabeth L. Turner, PhD,**

Department of Biostatistics and Bioinformatics, Duke University, Durham, NC

Duke Global Health Institute, Duke University, Durham, North Carolina, USA

**John A. Gallis, MSc,**

Department of Biostatistics and Bioinformatics, Duke University, Durham, NC

Duke Global Health Institute, Duke University, Durham, North Carolina, USA

**Fan Li, MS,**

Department of Biostatistics and Bioinformatics, Duke University, Durham, North Carolina, USA

**Melanie Prague, PhD,** and

Department of Biostatistics, Harvard T. H. Chan School of Public Health, Harvard University, Boston, Massachusetts, USA

Inria, project team SISTM, Bordeaux, France

**David M. Murray, PhD**

Office of Disease Prevention, Division of Program Coordination and Strategic Planning, Office of the Director, National Institutes of Health, Rockville, Maryland, USA

### Abstract

In 2004, Murray et al. published a review of methodological developments in both the design and analysis of group-randomized trials (GRTs). Over the last 13 years, there have been many developments in both areas. The goal of the current paper is to review developments in analysis, with a companion paper to focus on developments in design. As a pair, these papers update the 2004 review. This analysis paper includes developments in topics included in the earlier review, such as methods for parallel-arm GRTs, inference for conditional and marginal effects, and new topics including methods to account for multiple levels of clustering and alternative estimation methods such as augmented GEE, targeted maximum likelihood and quadratic inference functions.

---

**CORRESPONDING AUTHOR CONTACT INFORMATION** Elizabeth L. Turner, Ph.D., Assistant Professor, Department of Biostatistics and Bioinformatics and Duke Global Health Institute, Duke School of Medicine, Duke University, 2424 Erwin Road, Durham, NC 27710, USA. Phone: 919-681-6226. liz.turner@duke.edu.

#### ACCEPTANCE DATE

02/05/2017

#### CONTRIBUTORS

ELT and DMM initiated the project, developed the outline and topics to be covered, to which all authors agreed. ELT wrote much of the first draft, to which MP, JAG, FL and DMM contributed sections. All authors edited and reviewed the revised manuscript and approved it in its final version.

#### HUMAN PARTICIPANT PROTECTION

IRB approval was not needed.

We also examine developments in dealing with missing outcome data, including doubly robust approaches, software available for analysis, and analysis of alternative group designs (including stepped wedge GRTs, network-randomized trials, pseudo-cluster randomized trials and individually-randomized group treatment trials). These alternative designs, like the parallel-arm GRT, require clustering to be accounted for in both their design and analysis.

---

## INTRODUCTION

In a group-randomized trial (GRT), the unit of randomization is a group and outcome measurements are obtained on members of those groups.<sup>1</sup> Also called a cluster-randomized trial or community trial,<sup>2-5</sup> a GRT is the best comparative design available if the intervention operates at a group level, manipulates the physical or social environment, cannot be delivered to individual members of the group without substantial risk of contamination, or under other circumstances (e.g., a desire for herd immunity in studies of infectious disease).<sup>1-5</sup>

In GRTs, outcomes on members of the same group are likely to be more similar to each other than to outcomes on members from other groups.<sup>1</sup> Such clustering must be accounted for in the design to avoid an under-powered study and in the analysis to avoid underestimated standard errors and inflated type I error for the intervention effect.<sup>1-5</sup> For analysis, regression modeling approaches are generally preferred and most commonly used because of their ease of implementation.<sup>6</sup> Several textbooks now address these and other issues.<sup>1-5</sup>

In 2004, Murray et al.<sup>7</sup> published a review of methodological developments in both the design and analysis of GRTs. In the 13 years since, there have been many developments in both areas. The goal of the current paper is to focus on developments in analytic methods, including those relevant to designs described in a companion paper that focuses on developments in GRT design.<sup>8</sup> As a pair, these papers update the 2004 review. With both papers, we seek to provide a broad and comprehensive review to guide the reader to seek out appropriate materials for their own circumstances.

## DEVELOPMENTS IN THE ANALYSIS OF PARALLEL GROUP-RANDOMIZED TRIALS

### Methods for Superiority, Equivalence, and Non-Inferiority

In GRTs, superiority trials are more common than equivalence or non-inferiority trials: a PubMed search by one of the authors (DMM) of studies published in 2015 identified 562 superiority GRTs but only 1 equivalence GRT and 2 non-inferiority GRTs. Similarly, developments in the methods literature have focused on superiority GRTs, with developments for equivalence and non-inferiority GRTs limited to small sections in two of the more recent textbooks<sup>2,5</sup> and a review paper on sample size methods.<sup>9</sup> As a consequence, the current review paper focuses on superiority GRTs.

## Methods for Intention-To-Treat and Alternative Intervention Effects

In GRTs, protocol violations can lead to non-compliance at either the group- or member-level.<sup>5</sup> In order to minimize bias, intention-to-treat (ITT) principles are recommended at both levels rather than “on-treatment” and “per-protocol” analyses.<sup>2,4,5</sup> While group-level protocol violations are usually easy to identify, member-level compliance may be more difficult to ascertain in practice.<sup>2</sup> Jo et al. demonstrate that analyses which ignore compliance information could be underpowered to detect an ITT effect and propose a multilevel model combined with a mixture model.<sup>10</sup> Implications of group-level non-compliance can be considerable in GRTs, given the small number of groups that are randomized in many GRTs.

## Methods Based on the Randomization Scheme

Matching or stratification in the design has been recommended for some time as a way to ensure baseline balance on important potential confounders,<sup>1</sup> with constrained randomization more recently developed.<sup>11</sup> Recent reports suggest that most GRTs follow this advice.<sup>12–15</sup> Matching and stratification in the design can be ignored in the analysis of intervention effects, without harm to the type I error rate, and often the saved degrees of freedom will improve power.<sup>16,17</sup> Recently, Donner et al. reported that ignoring matching can adversely affect other analyses, such as analyses that examine the relationship between a risk factor and an outcome;<sup>18</sup> for this reason, investigators considering pair-matching should consider small strata instead (e.g., strata of 4). Li et al.<sup>19</sup> compared model-based and permutation methods in the context of constrained randomization adjusting for group-level covariates. They found that both the adjusted F-test and permutation test maintained the nominal size and had improved power under constrained randomization compared to simple randomization.

## Model-Based Methods

Model-based methods can be broadly classified according to the interpretation of the model parameters. Conditional model parameters are typically estimated using mixed-effects regression via maximum likelihood estimation (MLE) and are referred to as cluster-specific effects (or as subject-specific effects in the longitudinal analysis literature). Effects are conditional on the random effects used to account for clustering and on other covariates included in the analysis. Conditional models are often recommended for studies focused on change within members or on mediation analyses.<sup>7</sup> Parameters of marginal models are usually estimated using generalized estimating equations (GEE).<sup>20,21</sup> They define the marginal expectation of the dependent variable as a function of the independent variables and assume that the variance is a function of the mean; they separately specify a working correlation structure for observations made on members of the same group. Marginal models are often preferred for analyses of population-level effects because the intervention effect coefficient is interpreted as a population-averaged effect. In practice, marginal models are less frequently used than conditional models.<sup>6</sup>

Marginal and conditional intervention effects are equal for identity and log links<sup>22</sup> and the distinction between them is only important for link functions such as the logit for binary outcomes. Although some authors have advocated for the log instead of logit link for binary

outcomes,<sup>23</sup> this approach is not widely used, possibly because of model convergence problems for some data.<sup>24,25</sup> Alternatively, a modified Poisson approach with log-link and robust standard errors could be used in the GEE framework,<sup>26</sup> since it does not suffer from the same convergence problems as the binomial model with log link,<sup>27</sup> but it may be less common because of the familiarity of logistic regression among epidemiologists and biostatisticians.

In practice, the question about which of conditional or marginal effects are desired depends on the research question. It is essential to understand the underlying assumptions of each method: conditional models rely on correct specification of untestable aspects of the data distribution, while marginal models rely on a correct definition of the population of interest, which can make it difficult to generalize results to other populations.<sup>28</sup> We address each of the two approaches in more detail below.

**Conditional Approaches**—If the mixed effects model used to estimate conditional effects is misspecified, the estimates are difficult to interpret and, even if regression diagnostics can help,<sup>29</sup> standard errors (SEs) are not robust. Fortunately, Murray et al.<sup>30</sup> and Fu<sup>31</sup> have shown that mixed models are robust to substantial violation of the normality assumptions for member- and group-level errors, so long as balance is maintained at the group level. Parameter estimation by restricted maximum likelihood estimation (REML) is preferred to MLE when few groups are available.<sup>32–34</sup> For binary outcomes, alternative methods for specifying the test degrees of freedom have been examined in small sample GRTs and the between-within method is recommended.<sup>32,35</sup>

**Multiple Levels of Clustering in Conditional Models**—GRTs may involve multiple levels of clustering due to repeated measures on individuals or groups or additional hierarchical levels in the design. Murray<sup>1</sup> distinguished between mixed-effects models based on the number of measurements included in the analysis and recommended mixed-effects analysis of variance (ANOVA) or covariance (ANCOVA), or mixed-effects repeated measures ANOVA/ANCOVA, for analyses involving 1 or 2 measurements per person or per group; those models can account for all sources of random variation in such data if they are properly specified.<sup>36</sup> However, that is not the case in analyses involving 3 or more measurements per person or per group, where the sources of random variation may be different; instead, such analyses require a random coefficients model in which random trends and intercepts are calculated for each member (in cohort GRT designs) and group (in cohort and cross-sectional GRT designs), average trends and intercepts are calculated for each study arm, and the intervention effect is the net difference in the average study-arm trends.<sup>36</sup> Trends are often estimated as linear slopes, but can take another form.

**Variable Group Size in Conditional Models**—Johnson et al. focused on the analysis of Gaussian outcomes from GRTs with variable group size.<sup>37</sup> They compared ten model-based approaches and found that a one-stage mixed model with Kenward-Roger<sup>32</sup> degrees of freedom and unconstrained variance components performed well for GRTs with 14 or more groups per study arm. A two-stage model weighted by the inverse of the estimated theoretical variance of the group means and with unconstrained variance components

performed well for GRTs with 6 or more groups per study arm. A number of other models resulted in an inflated type I error rate when there was substantial variability in group size.

**Marginal Approaches**—When the GEE approach is used to estimate marginal effects, unbiased intervention effects can be estimated even if the working correlation structure is incorrect (e.g. using robust SEs via the sandwich estimator), although precision is increased if the working matrix is correct. Where degrees of freedom are limited for the test of interest, as often happens in GRTs, SE estimation is often biased downward and no method corrects for it in all cases, although several have been proposed.<sup>38–44</sup>

**Multiple Levels of Clustering in Marginal Models**—While multilevel clustering is easy to account for in mixed-effects regression, there is less literature for the GEE approach. The alternating logistic regression approach<sup>45</sup> for binary and ordinal outcomes can be used to account for correlation due to repeated measures on individuals within groups and can be implemented within a GEE framework in both R (the `alr` package) and SAS (PROC GEE).<sup>46</sup> The second-order GEE approach which, in contrast to regular GEE, models the working correlation structure as a function of covariates, can be implemented in R (`geepack` in R<sup>47</sup>).<sup>48</sup> For more general working correlation matrices, the user typically needs to perform additional programming in order to provide the appropriate covariance matrix and convergence may not be achieved. In addition, although the intervention effect is unbiased when the marginal model is not correctly specified, the SEs estimated using GEE may be too small. To correct this, a robust sandwich estimator of the variance can be used but such an approach leads to loss of power.<sup>49</sup> Because of this accuracy-power trade-off, mixed-effects models may be a better option to deal with GRTs involving more than two levels, although the effects estimated in such models are conditional rather than marginal effects.

**Variable Group Size in Marginal Models**—Although GEE analysis can accommodate variable group size, informative group size can negatively impact efficiency. In this case, Williamson et al.<sup>50</sup> showed that GEE weighted by group size can correct bias in the estimated intervention effect. This approach is equivalent and less computationally demanding than within-cluster resampling.<sup>51</sup>

**Advanced GEE Approaches to Improve Efficiency**—For binary outcomes, GEE is more conservative (i.e. the intervention effect will be estimated closer to the null) than mixed-effects models.<sup>28,52</sup> Moreover, the SE of the estimated intervention effect is also typically larger when using GEE so that much recent effort has focused on efficient estimation. GEE is most efficient when the true correlation structure of the data is chosen as the working correlation structure. Hin et al. compared multiple selection criteria for the working correlation matrix.<sup>53</sup> An alternative approach is augmented GEE (AU-GEE), a method developed for independent data using a causal inference framework,<sup>54</sup> which has been extended to clustered data.<sup>55</sup> AU-GEE uses covariate information to improve efficiency in a two-stage approach that specifies a model for the potential outcomes under the treatment not received. AU-GEE is unbiased and robust to misspecification of the potential outcome model, though correct specification improves efficiency. As for the analysis of all trials, only baseline covariates should be included in AU-GEE for the analysis of GRT data because

adjustment for post-baseline covariates may lead to bias.<sup>56</sup> Alternative methods are available to account for post-baseline, time-varying confounding.<sup>57–59</sup>

**Alternatives to GEE**—The quadratic inference function (QIF) method is an alternative to GEE for the estimation of marginal effects. Song et al.<sup>60</sup> demonstrate that QIF has advantages over GEE: it is more efficient and more robust to outliers; it has a goodness-of-fit test of the marginal mean model and permits straightforward extensions to model selection. In large samples, QIF is more efficient than GEE when the working correlation structure for the data is misspecified.<sup>61</sup> However, the SEs may be under-estimated for small and medium sample size or for variable group size.<sup>62</sup> More recent work by Westgate<sup>63,64</sup> provides improvements by using a bias-corrected sandwich covariance estimate and by simultaneously selecting the QIF or GEE while selecting the best working correlation structure.<sup>65</sup> Despite the many attractive properties of QIF, at this time there are few applications in public health.<sup>66–68</sup>

A second alternative estimation method is targeted maximum likelihood estimation (TMLE).<sup>69</sup> TMLE is a maximum likelihood-based G-computation estimator that targets the fit of the data-generating distribution to reduce bias in the parameter of interest. It is based on a machine learning approach that fluctuates an initial estimate of the conditional mean outcome and minimizes a loss function to provide an estimate of the parameter of interest.<sup>70</sup> The approach has been used in public health<sup>71,72</sup> and shows much promise for GRTs<sup>73,74</sup> because it can improve efficiency by simultaneously accounting for missing data and chance baseline covariate imbalance without committing to a specific functional form.<sup>75</sup>

### Permutation Methods

Permutation analysis was introduced for GRTs by Gail et al. for the COMMIT trial.<sup>76</sup> They found that the permutation test had nominal type I and II error rates across a variety of settings common to GRTs, when the member-level errors were Gaussian or binomial, even when very few heterogeneous groups were randomized to each study arm, and even when the ICC was large, so long as there was balance at the level of the group. Murray et al.<sup>30</sup> extended this work, showing that unadjusted permutation tests offer no more protection against confounding than unadjusted model-based tests, while the adjusted versions of both tests perform similarly. The permutation test was more powerful than the model-based test when the data were binomial and the ICC = 0.01. Fu<sup>31</sup> extended the work to heavy tailed and very skewed distributions and reported similar results.

Li et al. compared model-based and permutation methods in the context of constrained randomization adjusting for group-level covariates. They found that both the adjusted F-test and permutation test maintained the nominal size and had similar power, but cautioned that the randomization distribution must be calculated within the constrained randomization space to prevent inflating the type I error rate.<sup>19</sup>

## DEVELOPMENTS IN THE ANALYSIS OF ALTERNATIVES TO THE PARALLEL GRT

### Stepped Wedge GRT

Both between- and within-group information is available to estimate the intervention effect from a stepped wedge group randomized trial (SW-GRT).<sup>77,78</sup> However, because the control condition is typically observed earlier than the intervention condition, time is a potential confounder and should be accommodated in the analysis of SW-GRTs, typically by accounting for time as a predictor.<sup>79</sup> As for parallel GRTs, clustering by group must be accounted for, and longitudinal measures on individuals can be accommodated within either the mixed-effects or GEE framework, though more easily using mixed-effects models (see both *Multiple Levels of Clustering* sections). Conditional approaches are more commonly used in practice and reported on in the methods literature.<sup>79,80</sup> Several authors have highlighted other characteristics specific to SW-GRT including lagged intervention effects<sup>81</sup> and fidelity loss over time.<sup>79</sup>

### Network-Randomized GRT

Because the network properties of a network-randomized GRT are primarily used at the design stage,<sup>82</sup> and because they differ from regular GRTs only in the novel way in which groups are defined, the theory on the analysis of parallel-arm GRTs can be applied to parallel-arm network-randomized GRTs.<sup>83</sup> For example, in a ring trial of an Ebola vaccine,<sup>83</sup> in which a network was defined as all individuals who had regular physical contact with the incident (index) case of Ebola and in which all contacts received the vaccine (placebo or active), standard GRT methods were used. For network-randomized GRTs in which the intervention is not directly administered to all individuals and in which it is expected that the intervention spreads over the network (e.g. the snowball trials of a HIV prevention intervention for drug users<sup>84</sup> or a microfinance intervention<sup>85</sup>), methods<sup>86,87</sup> are available to estimate both the direct and indirect effects of the intervention. When network information is available and the outcome of interest is known to be a disseminated process, adjusting for network features such as information on the location of each individual within the network (i.e. group) can improve both the efficiency and power of the analysis.<sup>88</sup>

### Pseudo-Cluster Randomized Trial

Teerenstra et al.<sup>89</sup> compared analytic methods for continuous outcomes in pseudo-cluster randomized trials (PCRT) and Campbell and Walters discussed principles in their recent textbook.<sup>5</sup> Clustering by the unit of randomization at the first stage (e.g. provider) must be accounted for in both the design and analysis of PCRT. No explicit sample size or analytic methods are known to be available for non-continuous outcomes.

### Individually Randomized Group Treatment Trial

Baldwin et al. compared four analytic models for IRGTs and three methods for calculating degrees of freedom.<sup>90</sup> A multilevel model adapted to reflect clustering in only one study arm, combined with either Satterthwaite<sup>91</sup> or Kenward-Roger<sup>32</sup> degrees of freedom, provided better type I error control, better efficiency, and less bias, even with

heteroscedasticity at the member level. This finding is consistent with earlier reports by Pals et al.<sup>92</sup> and Roberts et al.<sup>93</sup> More recently, Roberts & Walwyn<sup>94</sup> and Andridge et al.<sup>95</sup> considered the circumstance in which members are associated with more than one small group or change agent. Both found that ignoring membership in multiple groups further inflates the type I error rate. Roberts & Walwyn reported that multiple member multilevel models maintained the nominal type I error rate; they also provide sample size and power formulae.<sup>94</sup>

## DEVELOPMENTS TO ADDRESS DATA CHALLENGES

### Missing Outcome Data

Two recent reviews<sup>6,96</sup> indicate that missing outcome data is common in GRTs, though investigators frequently analyze only available data without accounting for the missing data pattern. When the covariate-dependent missingness (CDM) assumption is plausible, both mixed effects and GEE models provide unbiased estimates of the intervention effect when the CDM covariates are included in an analysis of all available data.<sup>97,98</sup> AU-GEE also can provide unbiased effects by including all CDM covariates in the augmentation component<sup>55</sup> and has the advantage that all estimates can still be interpreted as marginal effects. Other two-stage approaches such as multiple imputation (MI) or inverse probability weighting (IPW) can provide unbiased intervention effects under certain conditions for more general missing at random (MAR) patterns and may provide increased precision compared to covariate-adjusted conditional or marginal models for CDM.<sup>97,99</sup> Although there is less literature on how to deal with missing not-at-random (MNAR) data,<sup>100</sup> sensitivity analyses are recommended.<sup>101</sup> A recent review showed that very few GRTs performed any sensitivity analyses for their missing data assumptions.<sup>6</sup>

To avoid possible type I error, MI should account for the clustered data structure.<sup>102,103</sup> Fixed group effects should not be used due to reduced power.<sup>104</sup> For binary outcomes, Ma et al.<sup>105</sup> and Caille et al.<sup>106</sup> show that the preferred MI method depends on the number of groups and the design effect, and note that bias may arise for some approaches even for CDM missingness. Using group-specific mean imputation may be adequate for continuous outcomes.<sup>98,102</sup> Hossain et al.<sup>98</sup> show that if the missing data mechanism has an interaction between a covariate predictive of the outcome and study arm, the imputation strategy must account for this interaction to be unbiased.

Whereas MI requires specifying the distribution of the missing data conditional on covariates, IPW requires specifying the probability of being missing depending on covariates. Theoretically, both approaches can be used for any type of outcome and for both CDM and more general forms of MAR mechanisms.<sup>99</sup> While IPW requires an additional assumption of positivity (all participants have a non-zero probability of being observed), it may be viewed as easier to define, particularly in the presence of non-intermittent missingness.<sup>107</sup> Importantly, and as for MI, if the missing data mechanism has an interaction between a covariate predictive of the outcome and study arm, the weights must be generated by accounting for this interaction in order to be unbiased.<sup>108</sup> Prague et al.<sup>109,110</sup> developed a doubly robust estimator in the context of IPW, which provides an unbiased estimate if either the marginal mean model or the missing data model is correctly specified. They



demonstrated that a doubly-robust augmented GEE approach can simultaneously account for both CDM and baseline covariate imbalance in GRTs when the parameter of interest is a marginal effect. Combining MI and IPW is a promising new approach which may have superior performance to IPW or MI alone when there are missing covariates in addition to missing outcomes.<sup>111</sup>

### Baseline Imbalance of Covariates

While design strategies such as restricted randomization<sup>8</sup> can help to achieve baseline covariate balance, they may not be easy to implement (e.g. if group characteristics are unknown in advance) and chance imbalance may arise regardless. In this case, some form of model-based covariate adjustment could be used such as standard multivariate regression for conditional models or AU-GEE for marginal models.<sup>55</sup> The advantage of AU-GEE in this case is that it is doubly robust in that the consistency of intervention effect estimate requires correct specification of either the marginal mean structure or the treatment model, and it separates covariate adjustment from intervention effect estimation thereby reducing the risk of choosing the adjustment models to obtain the most significant results. The standard multivariate regression adjustment approach does not enjoy either of these benefits.

Alternatively, Hansen and Bowers<sup>112</sup> proposed a balancing criterion and studied its randomization distribution in order to simultaneously test for balance of multiple covariates in both RCTs and GRTs. Leyrat et al.<sup>113</sup> suggested to use the c-statistic of the propensity score model to measure covariate balance at the individual level. Leon et al.<sup>114</sup> recommended propensity score matching to correct for baseline imbalance; in a simulation study, they report a median 90% reduction in bias. Nevertheless, the Consolidated Standards for Reporting of Trials (CONSORT)<sup>115</sup> recommends that the adjustment covariates be specified a priori for primary analyses so that secondary analyses could test sensitivity of the primary findings to adjustment for covariates identified post hoc.

### Software

Table 1 identifies three software programs that can be used to analyze data from GRTs. The table is organized around topics considered in the current paper. While none of the three software programs can readily implement both QIF and tMLE for GRTs, the R program offers the most ready-to-use functionality given its broad applicability to the methods cited in the current paper.

## REPORTING OF RESULTS

The CONSORT guidelines for individually randomized trials were extended to GRTs in 2004<sup>115</sup> and most journals now require authors to conform to these guidelines. Based on a review of 300 GRTs published between 2000–2008, Ivers et al. reported that 60% and 70% accounted for clustering in the sample size calculation and in the analysis, respectively, 56% used restricted randomization, and most (86%) allocated more than 4 groups per arm.<sup>14</sup> A more recent review of 86 trials published in 2013–2014 showed that 77% and 78% accounted for clustering in the sample size calculation and in the analysis, respectively, and that 51% used some form of restricted randomization.<sup>15</sup>

Given concerns about the ethical conduct of GRTs,<sup>116,117</sup> recent reports on conduct and reporting have focused on the ethics of GRTs. For example, Sim and Dawson discuss the challenges associated with obtaining informed consent in GRTs.<sup>118</sup> The Ottawa Statement on the ethical design and conduct of GRTs was published in 2012<sup>119</sup> with a reevaluation in 2015.<sup>120</sup>

## DISCUSSION

In this review, we have summarized many of the most important advances in the analysis of GRTs during the 13 years since the publication of the earlier review by Murray et al.<sup>7</sup> Many of these developments have focused on developments in marginal model parameter estimation (e.g. augmented GEE, QIF and tMLE) and missing data methods. Some topics that space limitations have prevented include review of recent developments in survival outcomes,<sup>2,121–125</sup> measurement bias,<sup>126,127</sup> validity,<sup>128,129</sup> Bayesian methods,<sup>4,130–132</sup> cost-effectiveness analyses<sup>4,133–136</sup> and mediation analyses to uncover mechanisms of action.<sup>137–140</sup>

Through this review, we have sought to ensure that the reader is reminded of the value of well-thought out analysis of GRTs and of keeping up to date with the many recent developments in this area. Pairing this knowledge with our companion review of developments in the design of GRTs,<sup>8</sup> we hope that our review leads to continued improvements in the design and analysis of GRTs.

## Acknowledgments

This work was partly funded by the following National Institutes of Health grants: R01 HD075875, R37 AI51164, R01 AI110478 and K01 MH104310. The authors would like to thank the two anonymous reviewers whose comments greatly helped improve the final version of this manuscript.

## APPENDIX: GLOSSARY

### Augmented GEE

“Augmenting the standard GEE with a function of baseline covariates.”<sup>55</sup> These methods adapt semiparametric theory developed by Robins<sup>141</sup> and Robins, Rotnitzky, and Zhao<sup>142</sup> for observational studies with time-varying exposures and missing data problems, respectively. They consist of leveraging the estimating equation by a predictor function for counterfactual outcomes under the intervention not received by the group/cluster considered missing.<sup>55</sup>

### Baseline covariate balance

The group-level and individual-level covariate distributions are similar in all study arms.<sup>11</sup>

### Choice of balancing criterion

Li et al. describe several balancing criteria to assess how well a GRT is balanced across covariates. These include the “best balance” (BB) metric of de Hoop et al.<sup>143</sup> the balance criterion (B) of Raab and Butcher,<sup>11</sup> and the total balance score introduced by Li et al.<sup>19</sup>

### Coefficient of variation

A measure of between-group variation, defined in Table 1 of our companion paper.<sup>8</sup>

**Cohort GRT design**

A cohort of individuals is enrolled at baseline and those same individuals are followed up over time.

**Constrained randomization**

Refers “to those designs that go beyond the basic design constraints to specify classes of randomization outcomes that satisfy certain balancing criteria, while retaining validity of the design.”<sup>144</sup>

**Cross-sectional GRT design**

A different set of individuals is obtained at each time point.

**Designed balance at the group level**

When there are equal numbers of groups randomized to each study arm.

**Intraclass correlation**

A measure of between-group variation, defined in Table 1 of our companion paper.<sup>8</sup>

**Covariate-dependent missingness (CDM) assumption**

The assumption that “missingness in outcomes depends on covariates measured at baseline, but not on the outcome itself.”<sup>98</sup>

**Doubly-robust augmented GEE approach**

Combining augmented GEE and IPW, a doubly-robust estimator is obtained, which provides an unbiased estimate if either the marginal mean model or the missing data model is correctly specified.<sup>109,110</sup>

**Equivalence**

Assessing whether the new intervention is equivalent to the comparison intervention.

**G-computation estimator**

A computational method to estimate causal effect in structural nested models. These models are designed to deal with confounding by variables affected by intervention.<sup>145</sup>

**Individually Randomized Group Treatment Trials**

“Studies that randomize individuals to study arms but deliver treatments in small groups or through a common change agent.”<sup>8,92</sup>

**Informative cluster size**

When the outcome measured is related to the size of the cluster.<sup>50</sup>

**Missing at Random (MAR) assumption**

Rubin’s (1976) definition is that “data are missing at random if for each possible value of the parameter  $\varphi$  [the parameter of the conditional distribution of the missing data indicator given the data], the conditional probability of the observed pattern of missing data, given the

missing data and the value of the observed data, is the same for all possible values of the missing data.”<sup>146</sup>

### **Network-Randomized GRT**

“The network-randomized GRT is a novel design that uses network information to address the challenge of potential contamination in GRTs of infectious diseases.”<sup>8,82,84,147</sup>

### **Non-inferiority**

When a trial is designed to show that the new intervention is not worse than the comparison intervention.

### **On treatment analyses**

When groups are analyzed “according to the intervention they actually received.”<sup>2</sup>

### **Per protocol analyses**

When groups “not receiving the correct intervention are excluded.”<sup>2</sup>

### **Pseudo-cluster randomized trial**

Intervention is allocated to individuals in a two-stage process. “In the first stage, providers are randomized to a patient allocation-mix.... In the second stage, patients recruited to the PCRT are individually randomized to intervention or control according to the allocation probability of their provider.”<sup>8</sup>

### **Stepped Wedge GRT**

“A one-directional crossover GRT in which time is divided into intervals and in which all groups eventually receive the intervention.”<sup>8,78</sup>

### **Superiority**

When a trial is designed to establish whether a new intervention is superior to the comparison intervention (e.g., another drug, a placebo, enhanced usual care). However, the statistical test is still two-sided, allowing for the possibility that the new intervention is actually worse than the comparison.

### **Within-cluster resampling**

Randomly sample one observation from each cluster, with replacement. Then analyze this resampled dataset. Repeat this process a large number of times. “The within-cluster resampling estimator is constructed as the average” of all of the resample-based estimates (see Hoffman et al.<sup>51</sup> pp. 1122-3).

## **References**

1. Murray, DM. Design and Analysis of Group-Randomized Trials. New York, NY: Oxford University Press; 1998.
2. Hayes, RJ., Moulton, LH. Cluster Randomised Trials. Boca Raton: CRC Press; 2009.
3. Donner, A., Klar, N. Design and Analysis of Cluster Randomization Trials in Health Research. London: Arnold; 2000.
4. Eldridge, S., Kerry, S. A Practical Guide to Cluster Randomised Trials in Health Services Research. Vol. 120. John Wiley & Sons; 2012.

5. Campbell, MJ., Walters, SJ. How to Design, Analyse and Report Cluster Randomised Trials in Medicine and Health Related Research. Chichester, West Sussex: John Wiley & Sons; 2014.
6. Fiero MH, Huang S, Oren E, Bell ML. Statistical analysis and handling of missing data in cluster randomized trials: a systematic review. *Trials*. 2016; 17(1):72. [PubMed: 26862034]
7. Murray DM, Varnell SP, Blitstein JL. Design and analysis of group-randomized trials: a review of recent methodological developments. *Am J Public Health*. 2004; 94(3):423–432. [PubMed: 14998806]
8. Turner EL, Li F, Gallis JA, Prague M, Murray DM. Review of Recent Methodological Developments in Group-Randomized Trials: Part 1 - Design. *Am J Public Health*. Submitted.
9. Rutterford C, Copas A, Eldridge S. Methods for sample size determination in cluster randomized trials. *Int J Epidemiol*. 2015; 44(3):1051–1067. [PubMed: 26174515]
10. Jo B, Asparouhov T, Muthén BO. Intention-to-treat analysis in cluster randomized trials with noncompliance. *Stat Med*. 2008; 27(27):5565. [PubMed: 18623608]
11. Raab GM, Butcher I. Balance in cluster randomized trials. *Stat Med*. 2001; 20(3):351–365. [PubMed: 11180306]
12. Varnell SP, Murray DM, Janega JB, Blitstein JL. Design and analysis of group-randomized trials: a review of recent practices. *Am J Public Health*. 2004; 94(3):393–399. [PubMed: 14998802]
13. Murray DM, Pals SP, Blitstein JL, Alfano CM, Lehman J. Design and analysis of group-randomized trials in cancer: a review of current practices. *J Natl Cancer Inst*. 2008; 100(7):483–491. [PubMed: 18364501]
14. Ivers NM, Halperin IJ, Barnsley J, et al. Allocation techniques for balance at baseline in cluster randomized trials: a methodological review. *Trials*. 2012; 13:120. [PubMed: 22853820]
15. Fiero M, Huang S, Bell ML. Statistical analysis and handling of missing data in cluster randomised trials: protocol for a systematic review. *BMJ Open*. 2015; 5(5):e007378.
16. Diehr P, Martin DC, Koepsell T, Cheadle A. Breaking the matches in a paired t-test for community interventions when the number of pairs is small. *Stat Med*. 1995; 14(13):1491–1504. [PubMed: 7481187]
17. Proschan MA. On the distribution of the unpaired t-statistic with paired data. *Stat Med*. 1996; 15(10):1059–1063. [PubMed: 8783442]
18. Donner A, Taljaard M, Klar N. The merits of breaking the matches: a cautionary tale. *Stat Med*. 2007; 26(9):2036–2051. [PubMed: 16927437]
19. Li F, Lokhnygina Y, Murray DM, Heagerty PJ, DeLong ER. An evaluation of constrained randomization for the design and analysis of group-randomized trials. *Stat Med*. 2015; 35(10):1565–1579. [PubMed: 26598212]
20. Liang KY, Zeger SL. Longitudinal data analysis using generalized linear models. *Biometrika*. 1986; 73(1):13–22.
21. Zeger SL, Liang K-Y. Longitudinal data analysis for discrete and continuous outcomes. *Biometrics*. 1986; 42(1):121–130. [PubMed: 3719049]
22. Ritz J, Spiegelman D. Equivalence of conditional and marginal regression models for clustered and longitudinal data. *Stat Methods Med Res*. 2004; 13(4):309–323.
23. Greenland S. Interpretation and choice of effect measures in epidemiologic analyses. *Am J Epidemiol*. 1987; 125(5):761–768. [PubMed: 3551588]
24. Blizzard L, Hosmer W. Parameter Estimation and Goodness-of-Fit in Log Binomial Regression. *Biom J*. 2006; 48(1):5–22. [PubMed: 16544809]
25. Williamson T, Eliasziw M, Fick GH. Log-binomial models: exploring failed convergence. *Emerging themes in epidemiology*. 2013; 10(1):1–10. [PubMed: 23375106]
26. Zou G, Donner A. Extension of the modified Poisson regression model to prospective studies with correlated binary data. *Stat Methods Med Res*. 2013; 22(6):661–670. [PubMed: 22072596]
27. Yelland LN, Salter AB, Ryan P. Performance of the modified Poisson regression approach for estimating relative risks from clustered prospective data. *Am J Epidemiol*. 2011; 174(8):984–992. [PubMed: 21841157]

28. Hubbard AE, Ahern J, Fleischer NL, et al. To GEE or not to GEE: comparing population average and mixed models for estimating the associations between neighborhood risk factors and health. *Epidemiology*. 2010; 21(4):467–474. [PubMed: 20220526]
29. Huang X. Diagnosis of Random-Effect Model Misspecification in Generalized Linear Mixed Models for Binary Response. *Biometrics*. 2009; 65(2):361–368. [PubMed: 18759837]
30. Murray DM, Hannan PJ, Varnell SP, McCowen RG, Baker WL, Blitstein JL. A comparison of permutation and mixed-model regression methods for the analysis of simulated data in the context of a group-randomized trial. *Stat Med*. 2006; 25(3):375–388. [PubMed: 16143991]
31. Fu, D. A comparison study of general linear mixed model and permutation tests in group-randomized trials under non-normal error distributions [Dissertation]. Memphis: Statistics, University of Memphis; 2006.
32. Kenward MG, Roger JH. Small sample inference for fixed effects from restricted maximum likelihood. *Biometrics*. 1997; 53(3):983–997. [PubMed: 9333350]
33. Localio AR, Berlin JA, Have TRT. Longitudinal and repeated cross-sectional cluster-randomization designs using mixed effects regression for binary outcomes: bias and coverage of frequentist and Bayesian methods. *Stat Med*. 2006; 25(16):2720–2736. [PubMed: 16345043]
34. Pinheiro, JC., Bates, DM. *Mixed-effects models in S and S-PLUS*. New York: Springer; 2000.
35. Li P, Redden DT. Comparing denominator degrees of freedom approximations for the generalized linear mixed model in analyzing binary outcome in small sample cluster-randomized trials. *BMC Med Res Methodol*. 2015; 15(1):38. [PubMed: 25899170]
36. Murray DM, Hannan PJ, Wolfinger RD, Baker WL, Dwyer JH. Analysis of data from group-randomized trials with repeat observations on the same groups. *Stat Med*. 1998; 17(14):1581–1600. [PubMed: 9699231]
37. Johnson JL, Kreidler SM, Catellier DJ, Murray DM, Muller KE, Glueck DH. Recommendations for choosing an analysis method that controls Type I error for unbalanced cluster sample designs with Gaussian outcomes. *Stat Med*. 2015; 34(27):3531–3545. [PubMed: 26089186]
38. McNeish D, Stapleton LM. Modeling clustered data with very few clusters. *Multivariate Behav Res*. 2016; 51(4):495–518. [PubMed: 27269278]
39. Li P, Redden DT. Small sample performance of bias-corrected sandwich estimators for cluster-randomized trials with binary outcomes. *Stat Med*. 2015; 34(2):281–296. [PubMed: 25345738]
40. Fay MP, Graubard BI. Small-Sample Adjustments for Wald-Type Tests Using Sandwich Estimators. *Biometrics*. 2001; 57(4):1198–1206. [PubMed: 11764261]
41. Mancl LA, DeRouen TA. A covariance estimator for GEE with improved small-sample properties. *Biometrics*. 2001; 57(1):126–134. [PubMed: 11252587]
42. Morel J, Bokossa M, Neerchal N. Small sample correction for the variance of GEE estimators. *Biom J*. 2003; 45(4):395–409.
43. Preisser JS, Lu B, Qaqish BF. Finite sample adjustments in estimating equations and covariance estimators for intracluster correlations. *Stat Med*. 2008; 27(27):5764–5785. [PubMed: 18680122]
44. Pan W, Wall MM. Small-sample adjustments in using the sandwich variance estimator in generalized estimating equations. *Stat Med*. 2002; 21(10):1429–1441. [PubMed: 12185894]
45. Carey V, Zeger SL, Diggle P. Modelling multivariate binary data with alternating logistic regressions. *Biometrika*. 1993; 80(3):517–526.
46. By K, Qaqish BF, Preisser JS, Perin J, Zink RC. ORTH: R and SAS software for regression models of correlated binary data based on orthogonalized residuals and alternating logistic regressions. *Comput Methods Programs Biomed*. 2014; 113(2):557–568. [PubMed: 24286728]
47. Halekoh U, Højsgaard S, Yan J. The R package geeppack for generalized estimating equations. *Journal of Statistical Software*. 2006; 15(2):1–11.
48. Crespi CM, Wong WK, Mishra SI. Using second-order generalized estimating equations to model heterogeneous intraclass correlation in cluster-randomized trials. *Stat Med*. 2009; 28(5):814–827. [PubMed: 19109804]
49. Teerenstra S, Lu B, Preisser JS, van Achterberg T, Borm GF. Sample size considerations for GEE analyses of three-level cluster randomized trials. *Biometrics*. 2010; 66(4):1230–1237. [PubMed: 20070297]

50. Williamson JM, Datta S, Satten GA. Marginal analyses of clustered data when cluster size is informative. *Biometrics*. 2003; 59(1):36–42. [PubMed: 12762439]
51. Hoffman EB, Sen PK, Weinberg CR. Within-cluster resampling. *Biometrika*. 2001; 88(4):1121–1134.
52. Neuhaus JM, Kalbfleisch JD, Hauck WW. A comparison of cluster-specific and population-averaged approaches for analyzing correlated binary data. *Int Stat Rev*. 1991; 59(1):25–35.
53. Hin L-Y, Carey VJ, Wang Y-G. Criteria for working–correlation–structure selection in GEE: Assessment via simulation. *Am Stat*. 2007; 61(4):360–364.
54. Tsiatis AA, Davidian M, Zhang M, Lu X. Covariate adjustment for two-sample treatment comparisons in randomized clinical trials: A principled yet flexible approach. *Stat Med*. 2008; 27(23):4658–4677. [PubMed: 17960577]
55. Stephens AJ, Tchetgen Tchetgen EJ, Gruttola VD. Augmented generalized estimating equations for improving efficiency and validity of estimation in cluster randomized trials by leveraging cluster-level and individual-level covariates. *Stat Med*. 2012; 31(10):915–930. [PubMed: 22359361]
56. Richiardi L, Bellocco R, Zugna D. Mediation analysis in epidemiology: methods, interpretation and bias. *Int J Epidemiol*. 2013; 42(5):1511–1519. [PubMed: 24019424]
57. Robins JM, Rotnitzky A, Zhao LP. Analysis of semiparametric regression models for repeated outcomes in the presence of missing data. *J Am Stat Assoc*. 1995; 90(429):106–121.
58. Robins JM, Greenland S, Hu F-C. Estimation of the causal effect of a time-varying exposure on the marginal mean of a repeated binary outcome. *J Am Stat Assoc*. 1999; 94(447):687–700.
59. Miglioretti DL, Heagerty PJ. Marginal modeling of multilevel binary data with time-varying covariates. *Biostatistics*. 2004; 5(3):381–398. [PubMed: 15208201]
60. Song PXX, Jiang Z, Park E, Qu A. Quadratic inference functions in marginal models for longitudinal data. *Stat Med*. 2009; 28(29):3683–3696. [PubMed: 19757486]
61. Khajeh-Kazemi R, Golestan B, Mohammad K, Mahmoudi M, Nedjat S, Pakravan M. Comparison of Generalized Estimating Equations and Quadratic Inference Functions in superior versus inferior Ahmed Glaucoma Valve implantation. *J Res Med Sci*. 2011; 16(3):235–244. [PubMed: 22091239]
62. Westgate PM, Braun TM. The effect of cluster size imbalance and covariates on the estimation performance of quadratic inference functions. *Stat Med*. 2012; 31(20):2209–2222. [PubMed: 22415948]
63. Westgate PM. A bias-corrected covariance estimate for improved inference with quadratic inference functions. *Stat Med*. 2012; 31(29):4003–4022. [PubMed: 22807168]
64. Westgate PM. A covariance correction that accounts for correlation estimation to improve finite-sample inference with generalized estimating equations: a study on its applicability with structured correlation matrices. *J Stat Comput Simul*. 2016; 86(10):1891–1900. [PubMed: 27818539]
65. Westgate PM. Criterion for the simultaneous selection of a working correlation structure and either generalized estimating equations or the quadratic inference function approach. *Biom J*. 2014; 56(3):461–476. [PubMed: 24431030]
66. Asgari F, Biglarian A, Seifi B, Bakhshi A, Miri HH, Bakhshi E. Using quadratic inference functions to determine the factors associated with obesity: findings from the STEPS Survey in Iran. *Ann Epidemiol*. 2013; 23(9):534–538. [PubMed: 23958406]
67. Bakhshi E, Etemad K, Seifi B, Mohammad K, Biglarian A, Koohpayehzadeh J. Changes in Obesity Odds Ratio among Iranian Adults, since 2000: Quadratic Inference Functions Method. *Comput Math Methods Med*. 2016; 2016:1–7.
68. Yang K, Tao L, Mahara G, et al. An association of platelet indices with blood pressure in Beijing adults: Applying quadratic inference function for a longitudinal study. *Medicine (Baltimore)*. 2016; 95(39):e4964. [PubMed: 27684843]
69. Van der Laan, MJ., Robins, JM. Unified methods for censored longitudinal data and causality. Springer Science & Business Media; 2003.
70. Gruber S, van der Laan MJ. A targeted maximum likelihood estimator of a causal effect on a bounded continuous outcome. *Int J Biostat*. 2010; 6(1):1–18.

71. Kotwani P, Balzer L, Kwarisiima D, et al. Evaluating linkage to care for hypertension after community-based screening in rural Uganda. *Trop Med Int Health*. 2014; 19(4):459–468. [PubMed: 24495307]
72. Ahern J, Karasek D, Luedtke AR, Bruckner TA, van der Laan MJ. Racial/ethnic differences in the role of childhood adversities for mental disorders among a nationally representative sample of adolescents. *Epidemiology*. 2016; 27(5):697–704. [PubMed: 27196805]
73. Balzer LB, Petersen ML, van der Laan MJ. Targeted estimation and inference for the sample average treatment effect in trials with and without pair-matching. *Stat Med*. 2016; 35(21):3717–3732. [PubMed: 27087478]
74. Schnitzer ME, van der Laan MJ, Moodie EE, Platt RW. Effect of breastfeeding on gastrointestinal infection in infants: a targeted maximum likelihood approach for clustered longitudinal data. *Ann Appl Stat*. 2014; 8(2):703–725. [PubMed: 25505499]
75. Van der Laan MJ, Polley EC, Hubbard AE. Super learner. *Stat Appl Genet Mol Biol*. 2007; 6(1)
76. Gail MH, Mark SD, Carroll RJ, Green SB, Pee D. On design considerations and randomization-based inference for community intervention trials. *Stat Med*. 1996; 15(11):1069–1092. [PubMed: 8804140]
77. Hemming K, Haines TP, Chilton PJ, Girling AJ, Lilford RJ. The stepped wedge cluster randomised trial: rationale, design, analysis, and reporting. *BMJ*. 2015; 350:h391. [PubMed: 25662947]
78. Spiegelman D. Evaluating public health interventions: 2. Stepping up to routine public health evaluation with the stepped wedge design. *Am J Public Health*. 2016; 106(3):453–457. [PubMed: 26885961]
79. Davey C, Hargreaves J, Thompson JA, et al. Analysis and reporting of stepped wedge randomised controlled trials: synthesis and critical appraisal of published studies, 2010 to 2014. *Trials*. 2015; 16(1):358. [PubMed: 26278667]
80. Mdege ND, Man M-S, Taylor CA, Torgerson DJ. Systematic review of stepped wedge cluster randomized trials shows that design is particularly used to evaluate interventions during routine implementation. *J Clin Epidemiol*. 2011; 64(9):936–948. [PubMed: 21411284]
81. Copas AJ, Lewis JJ, Thompson JA, Davey C, Baio G, Hargreaves JR. Designing a stepped wedge trial: three main designs, carry-over effects and randomisation approaches. *Trials*. 2015; 16(1):352. [PubMed: 26279154]
82. Harling G, Wang R, Onnela J, De Gruttola V. Leveraging contact network structure in the design of cluster randomized trials. *Clin Trials*. 2016 [Epub ahead of print].
83. Ebola ça Suffit Ring Vaccination Trial Consortium. The ring vaccination trial: a novel cluster randomised controlled trial design to evaluate vaccine efficacy and effectiveness during outbreaks, with special reference to Ebola. *BMJ*. 2015; 351:h3740. [PubMed: 26215666]
84. Latkin C, Donnell D, Liu TY, Davey-Rothwell M, Celentano D, Metzger D. The dynamic relationship between social norms and behaviors: the results of an HIV prevention network intervention for injection drug users. *Addiction*. 2013; 108(5):934–943. [PubMed: 23362861]
85. Banerjee A, Chandrasekhar AG, Duflo E, Jackson MO. The diffusion of microfinance. *Science*. 2013; 341(6144)
86. Ogburn EL, VanderWeele TJ. Causal diagrams for interference. *Stat Sci*. 2014; 29(4):559–578.
87. VanderWeele TJ, Tchetgen EJT, Halloran ME. Components of the indirect effect in vaccine trials: identification of contagion and infectiousness effects. *Epidemiology*. 2012; 23(5):751. [PubMed: 22828661]
88. Staples P, Prague M, Victor DG, Onnela J-P. Leveraging Contact Network Information in Clustered Randomized Trials of Infectious Processes. *arXiv preprint arXiv:1610.00039*. 2016
89. Teerenstra S, Moerbeek M, Melis RJ, Borm GF. A comparison of methods to analyse continuous data from pseudo cluster randomized trials. *Stat Med*. 2007; 26(22):4100–4115. [PubMed: 17328006]
90. Baldwin SA, Bauer DJ, Stice E, Rohde P. Evaluating models for partially clustered designs. *Psychological Methods*. 2011; 16(2):149–165. [PubMed: 21517179]
91. Satterthwaite FE. An approximate distribution of estimates of variance components. *Biometrics*. 1946; 2(6):110–114. [PubMed: 20287815]



92. Pals SP, Murray DM, Alfano CM, Shadish WR, Hannan PJ, Baker WL. Individually randomized group treatment trials: a critical appraisal of frequently used design and analytic approaches. *Am J Public Health*. 2008; 98(8):1418–1424. [PubMed: 18556603]
93. Roberts C, Roberts SA. Design and analysis of clinical trials with clustering effects due to treatment. *Clin Trials*. 2005; 2(2):152–162. [PubMed: 16279137]
94. Roberts C, Walwyn R. Design and analysis of non-pharmacological treatment trials with multiple therapists per patient. *Stat Med*. 2013; 32(1):81–98. [PubMed: 22865729]
95. Andridge RR, Shoben AB, Muller KE, Murray DM. Analytic methods for individually randomized group treatment trials and group-randomized trials when subjects belong to multiple groups. *Stat Med*. 2014; 33(13):2178–2190. [PubMed: 24399701]
96. Díaz-Ordaz K, Kenward MG, Cohen A, Coleman CL, Eldridge S. Are missing data adequately handled in cluster randomised trials? A systematic review and guidelines. *Clin Trials*. 2014; 11(5): 590–600. [PubMed: 24902924]
97. DeSouza CM, Legedza AT, Sankoh AJ. An overview of practical approaches for handling missing data in clinical trials. *J Biopharm Stat*. 2009; 19(6):1055–1073. [PubMed: 20183464]
98. Hossain A, Diaz-Ordaz K, Bartlett JW. Missing continuous outcomes under covariate dependent missingness in cluster randomised trials. *Stat Methods Med Res*. 2016
99. Seaman SR, White IR. Review of inverse probability weighting for dealing with missing data. *Stat Methods Med Res*. 2013; 22(3):278–295. [PubMed: 21220355]
100. Vansteelandt S, Rotnitzky A, Robins J. Estimation of regression models for the mean of repeated outcomes under nonignorable nonmonotone nonresponse. *Biometrika*. 2007; 94(4):841–860. [PubMed: 27453583]
101. Thabane L, Mbuagbaw L, Zhang S, et al. A tutorial on sensitivity analyses in clinical trials: the what, why, when and how. *BMC Med Res Methodol*. 2013; 13(1):92. [PubMed: 23855337]
102. Taljaard M, Donner A, Klar N. Imputation strategies for missing continuous outcomes in cluster randomized trials. *Biom J*. 2008; 50(3):329–345. [PubMed: 18537126]
103. Ma J, Akhtar-Danesh N, Dolovich L, Thabane L. Imputation strategies for missing binary outcomes in cluster randomized trials. *BMC Med Res Methodol*. 2011; 11(1):18. [PubMed: 21324148]
104. Andridge RR. Quantifying the impact of fixed effects modeling of clusters in multiple imputation for cluster randomized trials. *Biom J*. 2011; 53(1):57–74. [PubMed: 21259309]
105. Ma J, Raina P, Beyene J, Thabane L. Comparing the performance of different multiple imputation strategies for missing binary outcomes in cluster randomized trials: a simulation study. *J Open Access Med Stat*. 2012; 2:93–103.
106. Caille A, Leyrat C, Giraudeau B. A comparison of imputation strategies in cluster randomized trials with missing binary outcomes. *Stat Methods Med Res*. 2016; 25(6):2650–2669. [PubMed: 24713160]
107. Seaman S, Galati J, Jackson D, Carlin J. What is meant by “missing at random”? *Stat Sci*. 2013; 28(2):257–268.
108. Belitser SV, Martens EP, Pestman WR, Groenwold RH, Boer A, Klungel OH. Measuring balance and model selection in propensity score methods. *Pharmacoepidemiol Drug Saf*. 2011; 20(11): 1115–1129. [PubMed: 21805529]
109. Prague, M., Wang, R., De Gruttola, V. Harvard University Biostatistics Working Paper Series. Harvard University; 2016. CRTgeeDR: An R Package for Doubly Robust Generalized Estimating Equations Estimations in Cluster Randomized Trials with Missing Data.
110. Prague M, Wang R, Stephens A, Tchetgen Tchetgen E, DeGruttola V. Accounting for interactions and complex inter-subject dependency in estimating treatment effect in cluster-randomized trials with missing outcomes. *Biometrics*. 2016; 72(4):1066–1077. [PubMed: 27060877]
111. Seaman SR, White IR, Copas AJ, Li L. Combining multiple imputation and inverse-probability weighting. *Biometrics*. 2012; 68(1):129–137. [PubMed: 22050039]
112. Hansen BB, Bowers J. Covariate Balance in Simple, Stratified and Clustered Comparative Studies. *Stat Sci*. 2008; 23(2):219–236.

113. Leyrat C, Caille A, Foucher Y, Giraudeau B. Propensity score to detect baseline imbalance in cluster randomized trials: the role of the c-statistic. *BMC Med Res Methodol*. 2016; 16(1):9. [PubMed: 26801083]
114. Leon AC, Demirtas H, Li C, Hedeker D. Subject-level matching for imbalance in cluster randomized trials with a small number of clusters. *Pharm Stat*. 2013; 12(5):268–274. [PubMed: 23798334]
115. Campbell MK, Elbourne DR, Altman DG. CONSORT statement: extension to cluster randomised trials. *Br Med J*. 2004; 328(7441):702–708. [PubMed: 15031246]
116. Hutton JL. Are distinctive ethical principles required for cluster randomized controlled trials? *Stat Med*. 2001; 20(3):473–488. [PubMed: 11180314]
117. Taljaard M, Chaudhry SH, Brehaut JC, et al. Survey of consent practices in cluster randomized trials: improvements are needed in ethical conduct and reporting. *Clin Trials*. 2014; 11(1):60–69. [PubMed: 24346609]
118. Sim J, Dawson A. Informed consent and cluster-randomized trials. *Am J Public Health*. 2012; 102(3):480–485. [PubMed: 22390511]
119. Weijer C, Grimshaw JM, Eccles MP, et al. The Ottawa statement on the ethical design and conduct of cluster randomized trials. *PLoS Med*. 2012; 9(11)
120. van der Graaf R, Koffijberg H, Grobbee DE, et al. The ethics of cluster-randomized trials requires further evaluation: a refinement of the Ottawa Statement. *J Clin Epidemiol*. 2015; 68(9):1108–1114. [PubMed: 25910909]
121. Zeng D, Lin D, Lin X. Semiparametric transformation models with random effects for clustered failure time data. *Stat Sin*. 2008; 18(1):355–377. [PubMed: 19809573]
122. Cai T, Cheng S, Wei L. Semiparametric mixed-effects models for clustered failure time data. *J Am Stat Assoc*. 2002; 97(458):514–522.
123. Zhong Y, Cook RJ. Sample size and robust marginal methods for cluster-randomized trials with censored event times. *Stat Med*. 2015; 34(6):901–923. [PubMed: 25522033]
124. Zhan Z, de Bock GH, Wiggers T, Heuvel E. The analysis of terminal endpoint events in stepped wedge designs. *Stat Med*. 2016; 35(24):4413–4426. [PubMed: 27311403]
125. Xu, Z. *Statistical Design and Survival Analysis in Cluster Randomized Trials [Dissertation]*. The University of Michigan; 2011.
126. Kramer MS, Martin RM, Sterne JA, Shapiro S, Dahhou M, Platt RW. The double jeopardy of clustered measurement and cluster randomisation. *BMJ*. 2009; 339
127. Cho S-J, Preacher KJ. Measurement Error Correction Formula for Cluster-Level Group Differences in Cluster Randomized and Observational Studies. *Educ Psychol Meas*. 2016; 76(5):771–786.
128. Eldridge S, Ashby D, Bennett C, Wakelin M, Feder G. Internal and external validity of cluster randomised trials: systematic review of recent trials. *BMJ*. 2008; 336(7649):876–880. [PubMed: 18364360]
129. Caille A, Kerry S, Tavernier E, Leyrat C, Eldridge S, Giraudeau B. Timeline cluster: a graphical tool to identify risk of bias in cluster randomised trials. *BMJ*. 2016; 354
130. Ma J, Thabane L, Kaczorowski J, et al. Comparison of Bayesian and classical methods in the analysis of cluster randomized controlled trials with a binary outcome: the Community Hypertension Assessment Trial (CHAT). *BMC Med Res Methodol*. 2009; 9(1):37. [PubMed: 19531226]
131. Grieve R, Nixon R, Thompson SG. Bayesian hierarchical models for cost-effectiveness analyses that use data from cluster randomized trials. *Med Decis Making*. 2010; 30(2):163–175. [PubMed: 19675321]
132. Clark AB, Bachmann MO. Bayesian methods of analysis for cluster randomized trials with count outcome data. *Stat Med*. 2010; 29(2):199–209. [PubMed: 19856321]
133. Gomes M, Ng ES-W, Grieve R, Nixon R, Carpenter J, Thompson SG. Developing appropriate methods for cost-effectiveness analysis of cluster randomized trials. *Med Decis Making*. 2012; 32(2):350–361. [PubMed: 22016450]
134. Díaz-Ordaz K, Kenward M, Gomes M, Grieve R. Multiple imputation methods for bivariate outcomes in cluster randomised trials. *Stat Med*. 2016; 35(20):3482–3496. [PubMed: 26990655]

135. Ng ES, Diaz-Ordaz K, Grieve R, Nixon RM, Thompson SG, Carpenter JR. Multilevel models for cost-effectiveness analyses that use cluster randomised trial data: an approach to model choice. *Stat Methods Med Res.* 2013; 25(5):2036–2052. [PubMed: 24346164]
136. Díaz-Ordaz K, Kenward MG, Grieve R. Handling missing values in cost effectiveness analyses that use data from cluster randomized trials. *J R Stat Soc Ser A Stat Soc.* 2014; 177(2):457–474.
137. Hox JJ, Moerbeek M, Kluytmans A, van de Schoot R. Analyzing indirect effects in cluster randomized trials. The effect of estimation method, number of groups and group sizes on accuracy and power. *Front Psychol.* 2014; 5:78. [PubMed: 24550881]
138. MacKinnon DP, Fairchild AJ, Fritz MS. Mediation analysis. *Annu Rev Psychol.* 2007; 58:593–614. [PubMed: 16968208]
139. Vanderweele TJ, Hong G, Jones SM, Brown JL. Mediation and spillover effects in group-randomized trials: a case study of the 4Rs educational intervention. *J Am Stat Assoc.* 2013; 108(502):469–482. [PubMed: 23997375]
140. VanderWeele TJ. A unification of mediation and interaction: a 4-way decomposition. *Epidemiology.* 2014; 25(5):749–761. [PubMed: 25000145]
141. Robins, JM. Marginal structural models versus structural nested models as tools for causal inference. In: Halloran, ME., Berry, DA., editors. *Statistical models in epidemiology, the environment and clinical trials.* New York: Springer; 1999. p. 95-134.
142. Robins JM, Rotnitzky A, Zhao LP. Estimation of regression coefficients when some regressors are not always observed. *J Am Stat Assoc.* 1994; 89(427):846–866.
143. de Hoop E, Teerenstra S, van Gaal BG, Moerbeek M, Borm GF. The "best balance" allocation led to optimal balance in cluster-controlled trials. *J Clin Epidemiol.* 2012; 65(2):132–137. [PubMed: 21840173]
144. Moulton LH. Covariate-based constrained randomization of group-randomized trials. *Clin Trials.* 2004; 1(3):297–305. [PubMed: 16279255]
145. Vansteelandt S, Joffe M. Structural nested models and g-estimation: The partially realized promise. *Stat Sci.* 2014; 29(4):707–731.
146. Rubin DB. Inference and missing data. *Biometrika.* 1976; 63(3):581–592.
147. Staples PC, Ogburn EL, Onnela J-P. Incorporating Contact Network Structure in Cluster Randomized Trials. *Sci Rep.* 2015; 5:17581. [PubMed: 26631604]

**Table 1**

Summary of known functions and procedures to analyze GRTs using methods described in the current review.

Method	Software		
	SAS	Stata	R
<b>Outcomes analysis of all available data</b>			
Mixed-effects models	PROC MIXED	mixed	lme4
	PROC NL MIXED	melogit	nlme
	PROC GLIMMIX	mepoisson	
Generalized estimating equations (GEE)	PROC GENMOD <sup>1</sup>	xtgee	geeglm/geeM
Targeted maximum likelihood (tmLE)	N/A	N/A	N/A <sup>2</sup>
Quadratic inference function (QIF)	%qif	N/A	qif <sup>3</sup>
Permutation tests	%ptest	N/A	N/A
<b>Accounting for missing outcomes</b>			
Multiple imputation for clustered data	%mmi_impute <sup>4</sup>	REALCOM Impute	pan
	%mmi_analyze	mi impute <sup>4</sup>	jomo <sup>5</sup>
Inverse probability weighting (IPW)	PROC GENMOD <sup>6</sup>	N/A <sup>7</sup>	CRTgeeDR
<b>Causal-inference based methods<sup>8</sup></b>			
Augmented GEE (AU-GEE)	N/A	N/A	CRTgeeDR
Doubly robust AU-GEE	N/A	N/A	CRTgeeDR

Footnotes:

- <sup>1</sup>. PROC GEE is another option, but is in experimental phase and has limited usefulness for GRTs over and above PROC GENMOD.
- <sup>2</sup>. In R, tmle is available for tmLE, but at the time of writing, does not allow for clustering.
- <sup>3</sup>. As of the writing, the authors have been unable to load the package and it only allows equal cluster size, but Westgate has modified the code for GRTs with variable cluster size in the appendix of his paper<sup>63</sup>
- <sup>4</sup>. Only useful for continuous outcomes.
- <sup>5</sup>. In R, mice is available for multiple imputation but at the time of writing, does not account for clustering.
- <sup>6</sup>. Cannot account for imprecision in the weights.
- <sup>7</sup>. xtgee cannot accommodate individual-level weights but only group-specific weights.
- <sup>8</sup>. Both of the listed methods are related: AU-GEE accounts for baseline covariate imbalance and doubly robust AU-GEE, an extension of AU-GEE, accounts for both baseline covariate imbalance and missing data.

N/A: not available at the time of writing.