**"The long run effects of labor migration on human capital formation in communities of origin"**

**Taryn Dinkelman and Martine Mariotti**

**ONLINE APPENDIX**

**Appendix 1: Data**

This appendix describes the main data sources used in the paper and the construction of main outcome and explanatory variables

1. *Education and demographic variables from 1998 Census*
   - We use 100% microdata from the 1998 Census. These data are available from the Malawi National Statistics Office and from IPUMSI (https://international.ipums.org/international/).
   - Variables include: total years of schooling attained for everyone in the data, current geographic location (region anddistrict of the individual), age and gender. We create additional education variables: whether someone has attained any primary schooling, whether someone has completed primary schooling, and whether an individual reports being bilingual or not

2. *Education and demographic variables from the 1977 Census*
   - We digitize aggregate data tables constructed from the 100% microdata of the 1977 Census, reported in *Malawi 1977 Population Census Final Report Volumes I and II*, Malawi National Statistics Office, Zomba
   - Data are available at national, region and district level, and sometimes at district, sex and five-year cohort level.
   - Variables we use include: total years of schooling attained by each gender-five year age group at district level, the share of each district-gender-five year age group cell that has ever been to primary school, and the cell counts for each district-gender-five year age cell.
   We also use data on the number of men reporting a return from working abroad by district and five year age group, since the prior 1966 Census, and the number of boys and girls aged 10 to 19 who are employed outside the home, employed in the home, and enrolled in school

3. *Historic variables from older Census data*
   - Aggregate tables presented at the district level are available from published reports for the 1931(*Report on the Census of 1931*, Nyasaland Protectorate), 1945 (*Report on the Census of 1945*, Nyasaland Protectorate) and 1966 (*Malawi 1966 Population Census Final Report*, Malawi

National Statistics Office, Zomba) Malawian Census. We digitized various tables from these reports and matched them to current definitions of district boundaries

- Variables include: the log of population density in 1931 and 1945, the share of youth who are literate (English and the vernacular) in 1945, the fraction of men employed in different sectors (farming/non-farming, working for wages/no wages, unemployed) in 1966, and the number of adult men who work abroad in 1966, reported at the district level

4. *Geographic variables*

- Altitude: we compute altitude for each point on the Malawian grid map using data from the national map seamless server (http://seamless.usgs.gov/index.php) and the Viewshed tool in ArcGIS. We aggregate these measures to district level.

- We define areas of high, medium or low malaria susceptibility based on standard measures of altitude: high malaria areas (altitude below 650m), medium malaria areas (altitudes between 650m and 1100m) and low malaria areas (altitudes over 1100m).

- We create a district boundary crosswalk that links districts over time (across Census waves) and across name changes. We assign variables measured in earlier years to later Census district boundaries in this way:
    - For districts that were eventually combined in later years, we add district level values together
    - For districts split apart in later years, we apportion district totals to split districts using the fraction of physical area that each split district accounts for within the total district.

- We identify which districts contain a large tea or tobacco plantation using information in Christiansen (1984). The FAO's crop suitability index measuring whether a district is highly suitable for tobacco or tea production significantly predicts this estate district indicator. Because of the coarseness of this measure, there is uncaptured variation within estate districts in the prevalence of estate lands out of total district agricultural lands.

5. *Administrative data*

- Figure 1 is constructed using the location of Wenela/TEBA recruiting stations in 1937. We collected and digitized this historical data using a variety of sources. The main source included "Correspondence from the Secretariat, Zomba, Nyasaland 1935 (Circular number 8 1935, S1/169/35). We verified these stations were still open in later years using information from later Provincial Administration Reports (Northern Province: 7[th] December 1961 Ref. No. O.3.37 and Commissioner for Labour Circular, 25[th] March 1957)
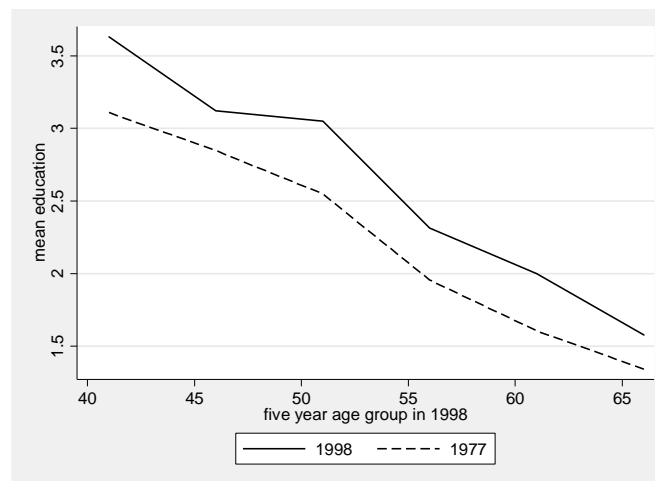
- Figure 2 is constructed using national labor migration totals from a variety of sources including: Chirwa (1991 for years 1950-1958); Lipton (1980: for years 1959-1994); Crush, Jeeves and Yudelman (1991: pp234-235) and various years of TEBA Annual Reports

**Appendix 2: Diagnosing selective attrition in the 1998 Census**

Because life expectancy in Malawi was only 46 years in the late 1990s (http://www.theglobaleconomy.com/Malawi/Life_expectancy/), we are concerned about mortality selection at ages over 40 affecting the composition of our sample. Specifically, we worry that this mortality selection may be differential across districts with and without recruiting stations. In this appendix, we use the 1977 and 1998 Census to diagnose this selective attrition and motivate our use of the 1977 Census data to construct estimates of educational attainment among the older cohorts in our analysis.[1]

Assuming that education is completed by age 20 in 1977, we can compare mean educational attainment for five year cohorts of those aged 20 to 44 in 1977 with the mean educational attainment of the five year age cohorts for those aged 41 to 65 in 1998.[2] At one extreme, if there is no mortality at all, mean education rates should match up for the same cohorts across Census waves. If there is any attrition (mortality) of those less educated between 1977 and 1998, then the mean education gap by cohort should be positive. Appendix 2 Figure 1 below shows just this.

**Appendix 2 Figure 1: Mean education by cohort in 1998 using 1977 and 1998 Census data**



Mean years of education (on the y axis) for each age group is higher in 1998 than for the corresponding age group in 1977, suggesting higher mortality rates among the less educated between 1977 and 1998.

---

[1] The question "What is the highest level of schooling you have attended?" is identical in the 1977 and 1998 Census' and the same coding system for different levels of education is used in both waves.

[2] Note that we cannot do this reliably for younger cohorts in 1977 (i.e for those 36-40 in 1998), since many of those under 20 are still in school.

Of greater concern for our identification strategy is differential mortality selection of the less educated in *Wenela* districts. The table below presents coefficients from a regression of the gap in mean completed years of education at district level (1998 levels minus 1977 levels) on dummies for each five year age cohort and interactions of each cohort dummy with number of recruiting stations in the district (the constant coefficient and coefficient on the number of recruiting stations are suppressed). All *Wenela* interaction terms are positive, and the interaction terms are jointly significantly different than zero. This means that the less-educated cohorts in *Wenela* regions are less likely to survive to 1998 than the less-educated cohorts in non-*Wenela* regions.

**Appendix 2 Table 1: Diagnosing selective attrition between 1977 and 1998**

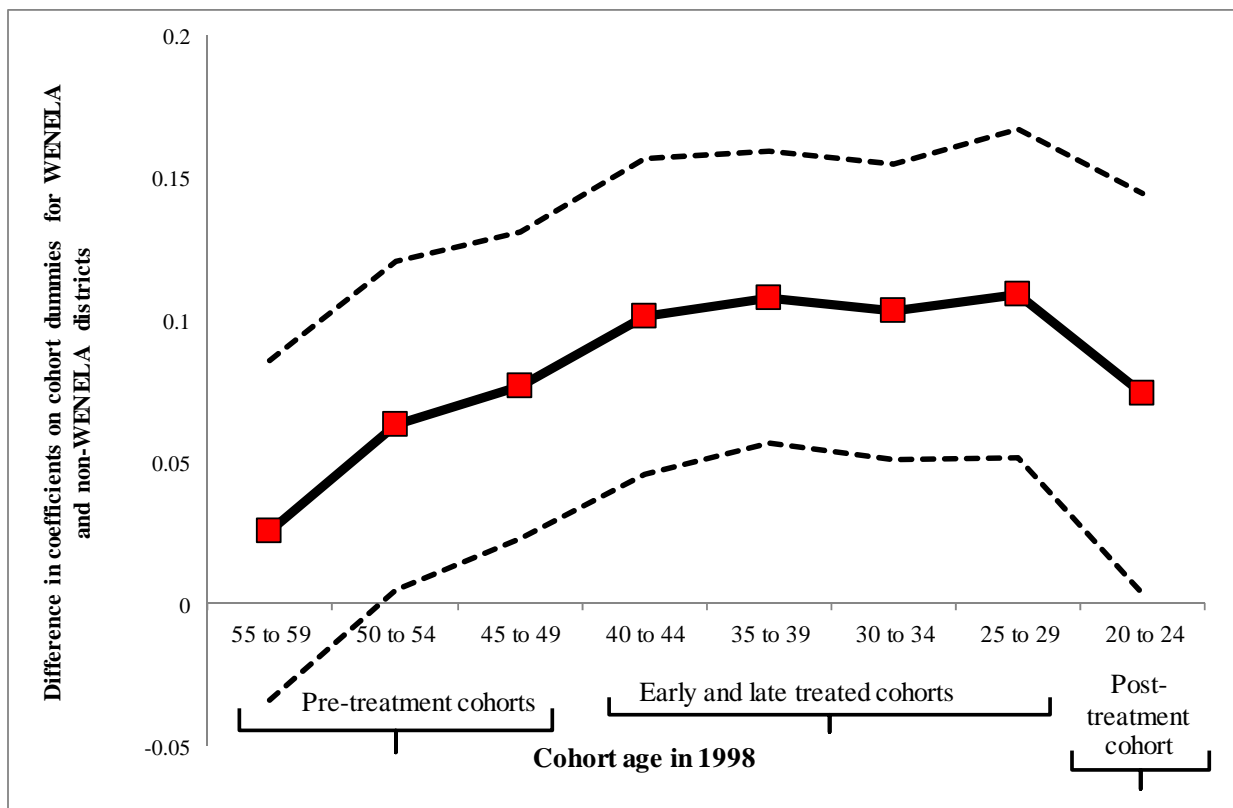| | |
|---|---|
| Age 40-44*Number of *Wenela* Stations | 0.0413 |
| | (0.0255) |
| Age 45-49*Number of *Wenela* Stations | 0.0288 |
| | (0.0212) |
| Age 50-54*Number of *Wenela* Stations | 0.0508*** |
| | (0.0192) |
| Age 55-59*Number of *Wenela* Stations | 0.0377** |
| | (0.0173) |
| Age 60-64*Number of *Wenela* Stations | 0.0313* |
| | (0.0183) |
| Age 40-44 | 0.348*** |
| | (0.0938) |
| Age 45-49 | 0.193** |
| | (0.0818) |
| Age 50-54 | 0.405*** |
| | (0.0878) |
| Age 55-59 | 0.322*** |
| | (0.0742) |
| Age 60-64 | 0.423*** |
| | (0.0821) |
| N | 135 |
| R2 | 0.654 |

N=189. Observation is the district-five year cohort. Outcome is the difference in mean level of education (1998 and 1977) measured for the district-five year cohort. Robust standard errors. F-statistic for interaction terms is 3.83 (*p* value is 0.003)

Together, these results suggest that using only the 1998 Census for our analysis would result in downwards-biased estimates of education gaps across *Wenela* and non-*Wenela* districts in the pre-treatment cohorts. That is, we would end up controlling for larger positive education gaps between the older cohorts in *Wenela* relative to non-*Wenela* districts, and would estimate smaller effects on educational attainment among the treatment cohorts in our difference-in-differences strategy.

We can illustrate the impact of ignoring the differential selective mortality across the *Wenela* and non-*Wenela* districts on our estimates. In Appendix 2 Figure 2, we present point estimates from the difference-

in-differences specification that generates our main result in Figure 4 of the paper. Instead of constructing the synthetic cohort of older pre-treatment age groups (ages 45 and over) from the 1977 Census as we do in Figure 4, we only use the 1998 Census to produce the graph below. The results in Appendix 2 Table 2 indicate that among the older age groups, less educated individuals have been differentially selected out in *Wenela* areas. Because we lose these people by 1998, the figure indicates an increase in the education gap between *Wenela* and non-*Wenela* cohorts even among the pre-treatment cohorts. This is the effect of differential selective attrition in *Wenela* districts.

**Appendix 2 Figure 2: Estimated differences in education by age group and Wenela status of district using 1998 cohorts only**



We still see the largest impacts on education gaps between *Wenela* and non-*Wenela* areas concentrated in the Early and Late treatment cohorts, as we do in our main results. However, the size of the effect is smaller because of the bias coming from pre-treatment differences in the education gap. The point estimates that correspond to the specification used to produce our main results in Table 5 but using only the 1998 Census in Appendix 2 Table 2. The difference in differences terms are still significantly different than zero for years of education, share with any primary schooling and two additional outcomes,

the share literate in English and the share bilingual (English and Chichewe), although smaller, as we might expect from the discussion in this appendix.

**Appendix 2 Table 2: Long run effects of labor migration shocks on education: Difference-in-differences results 1998 data only**

| | Total years of schooling attained | | Share with any primary schooling | | Share English literate | | Share bilingual | |
|---|---|---|---|---|---|---|---|---|
| | (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) |
| Early Treatment Cohorts*Num. Wenela stations ($\beta_1$) | 0.064** | 0.060*** | 0.004* | 0.004** | 0.006** | 0.005** | 0.006** | 0.005** |
| | (0.024) | (0.020) | (0.002) | (0.002) | (0.003) | (0.002) | (0.003) | (0.002) |
| Late Treatment Cohorts*Num. Wenela stations ($\beta_2$) | 0.066** | 0.059*** | 0.003 | 0.004 | 0.008** | 0.005*** | 0.008** | 0.005*** |
| | (0.028) | (0.018) | (0.002) | (0.003) | (0.004) | (0.002) | (0.004) | (0.002) |
| Post-Treatment Cohorts*Num. Wenela stations ($\beta_3$) | 0.034 | 0.026 | -0.002 | 0.000 | 0.004 | 0.001 | 0.004 | 0.001 |
| | (0.028) | (0.022) | (0.003) | (0.004) | (0.004) | (0.002) | (0.004) | (0.002) |
| Trend interactions | N | Y | N | Y | N | Y | N | Y |
| N | 432 | 432 | 432 | 432 | 432 | 432 | 432 | 432 |
| R2 | 0.96 | 0.97 | 0.95 | 0.95 | 0.95 | 0.96 | 0.95 | 0.96 |
| Mean of outcome variable | 3.66 | 3.66 | 0.60 | 0.60 | 0.36 | 0.36 | 0.36 | 0.36 |
| $p$ value of F test $H_0$: $\beta_1=\beta_2$ | 0.88 | 0.96 | 0.49 | 0.87 | 0.46 | 0.76 | 0.46 | 0.77 |
| $p$ value of F test $H_0$: $\beta_1=\beta_3$ | 0.17 | 0.15 | 0.04 | 0.17 | 0.32 | 0.11 | 0.32 | 0.11 |
| $p$ value of F test $H_0$: $\beta_2=\beta_3$ | 0.02 | 0.04 | 0.01 | 0.04 | 0.02 | 0.02 | 0.02 | 0.02 |

Robust standard errors clustered at the district level. Statistical significance at the 1, 5, and 10 percent levels is indicated by ***, **, and *, respectively, and evaluated relative to the small sample $t$ distribution to account for the small number of clusters. Unit of observation is the district-five year age cohort-gender cell. Vector of controls included in every regression: female dummy, cohort dummies, two region fixed effects, the log of district-level population density in 1931 and the share of literate youths in 1945, as well as district fixed effects. In the second specification for each outcome, trend terms are interacted with region fixed effects, baseline population density and baseline literacy rates. Number of Wenela stations in the district is a count variable. Sample includes adults ages 20 to 64 in the 1998 census.

**Appendix 3: Bounding results for composition effects from internal migration**

Internal migration poses one possible threat to the validity of our main results. Because neither the 1977 nor 1998 Census captures district of birth, we potentially mismeasure childhood exposure to Wenela recruiting stations among those people who move across districts after completing education, but before we see them in the relevant Census year. Internal migration flows are unlikely randomly allocated across districts. Without knowing more about differences in the magnitude and direction of migrant flows across districts, this possible misclassification of exposure to Wenela stations generates unpredictable biases in our estimates.

To see why, consider the following. Suppose all districts have the same average level of education before internal migration. If more educated adults move from non-Wenela to Wenela districts while less educated adults move in the opposite direction, this generates artificially positive differences in adult educational attainment across districts that we would ascribe to exposure to Wenela stations. As long as this sorting is constant over time, then the Wenela dummy in our regressions (as well as district fixed effects) controls for these observed differences in educational attainment driven by internal migration. If, however, internal migration flows differ by district as well as cohort then our results could be the result of complicated changes in the composition of population at the district level.

In the absence of individual level data on birth districts, we bound our effect sizes for possible composition changes induced by internal migration. We combine information on net migration rates from the 1977 Census with assumptions about possible levels of education of net migrants. First, we use 1977 Census data to construct the number of net migrants per person currently living in the district for each district in each five-year cohort and gender cell. We call this the net migration rate, or $NetMigRate_{asd}$. In our data, this number is always between -0.35 and 0.29.[1] We need to assume that this net migration rate is the same in 1977 and 1998, since the 1998 Census contains no information on district of birth. Second, we assume that all migrants – whether they show up as in- or outmigrants in a particular district – have

---

[1] Census 1977 counts the number of people in each district, cohort, and gender cell and enumerates how many of these individuals were born in each district. The net migration rate is computed as the difference between total in-migrants and total out-migrants divided by total current population in the district; it is the number of net migrants (in-migrants – out-migrants) per person living in the district. A 0.2 net migration rate means that for every person living in the district, there are 0.2 net in-migrants.

the same level of education and therefore we need only account for the potential education of net migrants, the difference between in- and outmigrants.[2]

We adjust our education variables ($\bar{Y}_{asd}$) measured at district, cohort, and gender level:

$$\bar{Y}_{asd}^{BOUND} = \frac{N_{asd}\bar{Y}_{asd} - NetMigrants_{asd}*\bar{Y}_{as}^m}{N_{asd} - NetMigrants_{asd}} \qquad (A.1)$$

$$= \frac{N_{asd}\bar{Y}_{asd} - NetMigRate_{asd}*N_{asd}*\bar{Y}_{as}^m}{N_{asd} - NetMigRate_{asd}*N_{asd}}$$

$$= \frac{\bar{Y}_{asd} - NetMigRate_{asd}*\bar{Y}_{as}^m}{1 - NetMigRate_{asd}}$$

where *BOUND = {upper, lower}*, $\bar{Y}_{asd}^{BOUND}$ represents the adjusted mean education outcome at district, cohort, and gender level, $N_{asd}$ is total population in a district-cohort-gender cell, $\bar{Y}_{as}^m$ is either the maximum or minimum value of the relevant education variable across all districts at cohort and gender level, and *NetMigrants$_{asd}$* is the total number of net migrants in a district-cohort-gender cell. *NetMigrants$_{asd}$* is estimated by multiplying the total population in that district-cohort-gender cell with the net migration rate (*NetMigRate$_{asd}$*) for that cell. Each component of (A.1) comes from the relevant Census wave, except for *NetMigRate$_{asd}$* which is computed using 1977 Census data and applied to both Census waves. We estimate the main regression specifications for our sample after creating these adjusted education variables, one set for each of the extreme values of $\bar{Y}_{as}^m$, or

$$\bar{Y}_{asd}^{lower} = \frac{\bar{Y}_{asd} - NetMigRate_{asd}*\bar{Y}_{as}^{max}}{1 - NetMigRate_{asd}} \text{ and}$$

$$\bar{Y}_{asd}^{upper} = \frac{\bar{Y}_{asd} - NetMigRate_{asd}*\bar{Y}_{as}^{min}}{1 - NetMigRate_{asd}}.$$

There are two notable features of equation (A.1). First, the adjustments we make for internal migration imply that $\bar{Y}_{asd}^{upper}$ and $\bar{Y}_{asd}^{lower}$ provide upper and lower bounds on mean education and average share of adults with any primary school across the entire sample. Second, despite these names, these adjustments do not imply that the difference-in-differences regressions using these new variables will produce estimates that contain our main education results. This is because in a closed system (i.e. the whole of Malawi) some districts are receiving districts (*NetMigRate$_{asd}$>0*) while others are sending districts

---

[2] For example: if there are 110 in-migrants and 100 out-migrants to a particular district, and in-migrants and out-migrants have the same levels of education, the only change in composition that occurs as a result of this net migration is due to the additional 10 people who migrated into the district.

($NetMigRate_{asd}<0$). In order for $\bar{Y}_{asd}^{upper}>\bar{Y}_{asd}$ or $\bar{Y}_{asd}^{lower}<\bar{Y}_{asd}$, the following equations should hold (note that in our sample, $1 - NetMigRate_{asd} > 0$ in all cases):

$$NetMigRate_{asd} * \left(\bar{Y}_{asd} - \bar{Y}_{as}^{min}\right) > 0 \qquad \text{(A.2)}$$

$$NetMigRate_{asd} * \left(\bar{Y}_{asd} - \bar{Y}_{as}^{max}\right) < 0 \qquad \text{(A.3)}$$

Since $\bar{Y}_{asd} \geq \bar{Y}_{as}^{min}$ and $\bar{Y}_{asd} \leq \bar{Y}_{as}^{max}$ in all districts, these equations are only satisfied for receiving districts that have $NetMigRate_{asd} > 0$. To see this, assume that we impute the minimum level of education for net migrants, Then, $\bar{Y}_{asd}^{upper} > \bar{Y}_{asd}$ is satisfied only in receiving districts because our adjustments take out the low levels of education of net in-migrants to create a higher adjusted mean education variable. For sending districts, where $NetMigRate_{asd} < 0$, the inequality in (A.2) is reversed and $\bar{Y}_{asd}^{upper} < \bar{Y}_{asd}$. Similarly, when we impute the maximum level of education for net migrants, equation (A.3) will only be satisfied in receiving districts; subtracting high levels of net in-migrant education generates $\bar{Y}_{asd}^{lower} < \bar{Y}_{asd}$. In sending districts, (A.3) is reversed, so $\bar{Y}_{asd}^{lower} > \bar{Y}_{asd}$.

Because we have both sending *and* receiving districts in our sample, and because rates of internal migration in 1977 are different across Wenela and non-Wenela districts (rates of in-migration are higher in Wenela districts, results not shown), our adjustments have different effects on the bounds values in specific Wenela and non-Wenela districts. More complicated patterns of net migration that vary across exposed and non-exposed cohorts and across Wenela and non-Wenela areas imply that adjustments for internal migration may generate in difference-in-differences estimates that do not bound our main result.[3] Nevertheless, it is still a useful exercise to check whether internal migration modelled in this way appears to confound our results.

Appendix Table 2 displays results from difference-in-differences regressions estimated using the adjusted education variables, first including all controls and district fixed effects, and then adding in trend interactions with region fixed effects, baseline literacy and baseline population density. We compare the coefficients in this table with the main estimates in Table 5.

First, assuming net migrants have the maximum level of schooling in the district-cohort-gender cell for a given Census year, the presence of anyone new in a receiving district raises mean education and their absence from a birth district artificially deflates that district's average education. Adjusting for these

---

[3] Crudely, if net migration rates are more likely positive in Wenela districts among exposed cohorts, we would be doing more "receiving district" adjustments in our core treatment groups and more "sending district" adjustments in our control groups.

educated net migrants, we still see large, positive impacts of exposure to treatment among exposed cohorts: those exposed during the labor expansion have 0.14 more years of education, while those exposed during the labor contraction have 0.2 more years of education. The education gap between Wenela and non-Wenela districts for the post-treatment group continues to be positive, at around 0.16 more years of education. Second, if we instead assume that net migrants are uneducated, removing them from our outcome measure in receiving districts and adding them back to sending districts reveals similar, large positive impacts of exposure to mine employment shocks. The difference-in-differences estimates in columns (3) and (4) imply that exposed cohorts in Wenela areas gained between 0.16 and 0.25 more years of education. Post-treatment groups continue to have about 0.22 more years of education in Wenela relative to non-Wenela areas, controlling for differences between these areas using the oldest pre-treatment cohorts. These bounds compare favorably to our main results in Table 5, 0.12 and 0.179 more years of education (Table 5, column 3).

Results are similar when we use the share with any primary school as outcome. In Table 5, directly exposed cohorts from districts with more Wenela stations are 1.1 to 2.6 percentage points more likely to have ever attended primary school. After adjusting for internal migration in Appendix Table 2 columns (5)-(8), these exposed cohorts from districts with Wenela stations are between 1.2 and 2.4 percentage points more likely to have ever been to primary school. All of our estimates are statistically different from zero at the 1, 5 or 10% level. Looking at the final three rows of the table, we see that we can strongly reject that the impacts of the migration shock on education are the same for the labor expansion and labor contraction cohorts, and we can reject that the impacts on primary school access for the labor contraction cohorts and the youngest post-treatment cohorts are the same. However, as in the case of our main results in Table 5, we cannot reject that the impacts for the post-treatment cohorts and the labor expansion cohorts are the same. We see the same inverted U-shaped pattern of coefficients in our bounded results as in the main results. The results of this bounding exercise suggest that selective internal migration and any resulting measurement error in $Wenela_d$ cannot account for our main effects.

**Appendix 3 Table 1: Long run effects of labor migration shocks on education: Bounds for internal migration**

| Assumptions about migrant education: | Total years of education | | | | Share with any primary school | | | |
|---|---|---|---|---|---|---|---|---|
| | Max. schooling | | Min. schooling | | Highest share with primary school | | Lowest share with primary school | |
| | (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) |
| Early Treatment Cohorts*Num. Wenela stations ($\beta_1$) | 0.142*** | 0.136*** | 0.163** | 0.164** | 0.014*** | 0.012*** | 0.015*** | 0.014*** |
| | (0.048) | (0.044) | (0.059) | (0.063) | (0.005) | (0.004) | (0.005) | (0.005) |
| Late Treatment Cohorts*Num. Wenela stations ($\beta_2$) | 0.202*** | 0.191*** | 0.251*** | 0.252** | 0.020*** | 0.016*** | 0.024*** | 0.021** |
| | (0.067) | (0.060) | (0.084) | (0.090) | (0.007) | (0.005) | (0.007) | (0.007) |
| Post-Treatment Cohorts*Num. Wenela stations ($\beta_3$) | 0.168** | 0.153** | 0.220*** | 0.222** | 0.015** | 0.011* | 0.019*** | 0.015* |
| | (0.062) | (0.062) | (0.076) | (0.088) | (0.005) | (0.006) | (0.006) | (0.007) |
| District FE | Y | Y | Y | Y | Y | Y | Y | Y |
| Region trends | N | Y | N | Y | N | Y | N | Y |
| N | 480 | 480 | 480 | 480 | 480 | 480 | 480 | 480 |
| R2 | 0.83 | 0.86 | 0.83 | 0.85 | 0.824 | 0.84 | 0.82 | 0.83 |
| Mean of outcome variable | 2.53 | 2.53 | 2.59 | 2.59 | 0.41 | 0.41 | 0.41 | 0.41 |
| $p$ value of F test H$_0$: $\beta_1=\beta_2$ | 0.01 | 0.00 | 0.00 | 0.01 | 0.01 | 0.02 | 0.00 | 0.02 |
| $p$ value of F test H$_0$: $\beta_1=\beta_3$ | 0.22 | 0.56 | 0.02 | 0.09 | 0.56 | 0.69 | 0.10 | 0.76 |
| $p$ value of F test H$_0$: $\beta_2=\beta_3$ | 0.01 | 0.09 | 0.04 | 0.12 | 0.01 | 0.07 | 0.01 | 0.04 |

Statistical significance at the 1, 5, and 10 percent levels is indicated by ***, **, and *, respectively, and evaluated relative to the small sample $t$ distribution to account for the small number of clusters. Robust standard errors clustered at the district level. Unit of observation is the district-5 year age group-gender cell. Vector of controls includes female, age group dummies, two region fixed effects, the log of district-level population density in 1931 and the share literate in 1931, and historical density, historical literacy and region fixed effects interacted with a trend term. Number of Wenela stations in the district is a count variable. Outcomes are our estimates of the bounds on education and share in primary school, after accounting for maximum and minimum possible values of each variable for the number of net migrants in each age-gender cell. Details of variable construction are explained in the text. Sample includes adults ages 20 to 44 in 1977 and 1998 census.

**Appendix 4: Constructing exact *p* values for Table 5 results using randomization inference**

Our empirical strategy exploits pre-existing spatial variation in migration costs across 24 districts within Malawi. The relatively small number of districts leads to a concern that standard inference procedures will over-reject the null hypothesis of zero impact of district-level exposure to the labor migration shock. We deal with this concern by presenting robust standard errors clustered at the district-level in Table 5 and obtain *p* values to indicate statistical significance using the small sample *t* distribution adjusted for the number of covariates that are constant within the cluster.

An alternative approach is to construct exact *p* values using randomization inference (Fisher, 1935; Rosenbaum 2002; see Cohen and Dupas 2010 for an example of how this is done in the context of a randomized controlled trial with 16 clusters and a single treatment). The idea behind this approach is as follows:

- We randomly assign the actual distribution of Wenela stations to districts and estimate the difference-in-difference models of Table 5 for this "false" assignment. The false allocation mimics the true distribution of stations

- Since there are over 1.3 million ways to allocate Wenela stations (the range of this variable is 0 to 10) to the 24 districts, we generate 1,000 different random assignments of stations to districts and estimate the difference-in-differences regression for each allocation.

- We compute the empirical distribution of *t* statistics for each of the three main parameters ($\beta_1$, $\beta_2$ and $\beta_3$) generated by these 1,000 false assignments

- We compare the actual *t* statistics from Table 5 to the empirical distribution of test statistics for each parameter and compute the probability of observing a *t* statistic in the tails of this distribution. The resulting *p* values, denoted randomization inference *p* values are presented in the table below.

In all cases, we can reject the null of zero impact at the 10% level. In most cases, for the estimates of $\beta_1$ and $\beta_2$, we can also reject the null of zero impact at the 5% level.

**References**

Cohen, Jessica and Pascaline Dupas, 2010 "Free distribution or cost-sharing? Evidence from a randomized malaria prevention experiment", *Quarterly Journal of Economics*, Vol. 125 (1): 1-45
Fisher, Ronald A. 1935 The Design of Experiments London: Oliver and Boyd.
Rosenbaum, Paul R. 2002 Observational Studies New York: Springer-Verlag.

**Appendix 4 Table 1: Randomization inference *p* values for difference-in-differences estimates of Table 5**

| *p* value for: | Outcome is: Years of schooling | | Outcome is: Share with any primary school | |
|---|---|---|---|---|
| $H_0: \beta_1=0$ | 0.001 | 0.001 | 0.001 | 0.016 |
| $H_0: \beta_2=0$ | 0.009 | 0.012 | 0.013 | 0.011 |
| $H_0: \beta_3=0$ | 0.013 | 0.031 | 0.009 | 0.084 |
| Other controls? | Y | Y | Y | Y |
| District FE? | Y | Y | Y | Y |
| Trend interactions | N | Y | N | Y |

**Appendix 5: Correlations between human capital formation and household well-being in Malawi**

Given concerns about the quality of education provided in poor countries, it is worthwhile asking whether there is any evidence that more years of schooling and more skills learned in school are valuable in Malawi. We use data from the 2004 Malawi Integrated Household Survey to estimate the relationship between total assets owned by the household (as a summary measure of household well-being) and human capital attainment of the head of the household. There are too few wage earners in the data to examine the relationship between individual education and wages.

We restrict the sample to all rural households in the survey where the head is between the ages of 26 and 60 in 2004, leaving us with 85% of the initial survey. We calculate total assets owned by the household (the range is 0 to 31; the mean is 6.9) and regress this index on four different measures of human capital of the household head: total years of schooling attained, whether the head has been to primary school, and self-reported literacy in English and Chichewa (the person reports that they can read a one page letter in the relevant language). Using the education of the head of household (rather than the maximum level of education in the household) is a meaningful way to characterize how much human capital the household has access to, since there are few three-generation households in the sample. Literacy reflects skills acquired in the lower levels of primary school: 88% of household heads with four years of completed education report being literate in Chichewa, and literacy rates increase from 47% to 72% between two and three years of total schooling.

Appendix 5 Table 1 presents the results of regressing our asset index on the four human capital measures for three specifications: one without controls, a second adding controls for age and gender of the household head, household size, historical district characteristics (literacy, population density) and region fixed effects, and a third that includes district fixed effects. Results indicate a robust, large and positive relationship between each human capital measure and the household asset index. In the first three columns of the table, we see that an additional year of schooling is correlated with about one third more total assets (columns 1-3), which is a 5-6% return per year of schooling. For a household head with any primary school, assets are 7% higher (columns 4 and 5). Returns to literacy are particularly large: having a literate household head raises total assets in the household by over two; a 30% gain. While usual concerns about selection and measurement error caution us against interpreting these point estimates as exactly causal, the strong positive relationship between human capital attained in childhood and measures

of household well-being in adulthood provide some evidence that education, and skills learned at school, are indeed valuable in Malawi.[1]

---

[1] Our results are consistent with findings in Chirwa and Matita (2009), who use Mincerian regressions to show that among wage earners working in urban areas of Malawi, the return to completed primary education is 5 percent.

**Appendix 5 Table 1: Correlation between household asset index and human capital of household head**

| | (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) | (9) | (10) | (11) | (12) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | *Outcome: Total assets owned by household [mean=6.9]* | | | | | | | | | | | |
| *Human capital of household head* | | | | | | | | | | | | |
| Years of education | 0.380*** | 0.362*** | 0.361*** | | | | | | | | | |
| | (0.023) | (0.023) | (0.019) | | | | | | | | | |
| Any primary school | | | | 0.253 | 0.471*** | 0.473*** | | | | | | |
| | | | | (0.165) | (0.148) | (0.146) | | | | | | |
| Literate in English | | | | | | | 2.531*** | 2.217*** | 2.185*** | | | |
| | | | | | | | (0.129) | (0.130) | (0.130) | | | |
| Literate in Chichewa | | | | | | | | | | 2.435*** | 2.061*** | 2.005*** |
| | | | | | | | | | | (0.148) | (0.166) | (0.147) |
| | | | | | | | | | | | | |
| N | 6,771 | 6,771 | 6,771 | 6,771 | 6,771 | 6,771 | 6,794 | 6,794 | 6,794 | 6,794 | 6,794 | 6,794 |
| Effect size: % of mean assets | 6% | 5% | 5% | 4% | 7% | 7% | 37% | 32% | 32% | 35% | 30% | 29% |
| | | | | | | | | | | | | |
| Additional controls? | N | Y | Y | N | Y | Y | N | Y | Y | N | Y | Y |
| District FE? | N | N | Y | N | N | Y | N | N | Y | N | N | Y |

Robust standard errors clustered at the district level. Statistical significance at the 1, 5, and 10 percent levels is indicated by ***, **, and *, respectively. Unit of observation is the household; household head's level of education is the main regressor in each specification. Sample includes all rural households where the household head is between the ages of 26 and 60 in 2004. Additional controls include age of household head, whether the head is female, household size, the log of district-level population density in 1931, the share of literate youths in the district in 1945, and region fixed effects. All regressions are weighted using household weights. Data are from the 2004/2005 Malawi Integrated Household Survey.