

Online Appendix to “Are Ideas Getting Harder to Find?”

Nicholas Bloom, Charles I. Jones, John Van Reenen, and Michael Webb

(Not for publication)

February 26, 2018

1. Overview

This document provides more details on the data used in our paper and the programs that can be used to replicate our results.

The programs/data for each portion of the paper are stored in a separate subdirectory. In general, the basic data are contained in spreadsheet files, and matlab programs are used to conduct the analysis. We use Matlab 2017a.

The following program can be run in the main directory where “IdeaPFPrograms.zip” is unzipped:

- **MasterIdeaPF.m:** Master program for generating all the results in the paper.

Note that you will need to edit this file to change to the main directory and to add the proper path to the “ChadMatlab” directory that is unzipped from IdeaPFPrograms.zip

2. Aggregate U.S. Evidence

The analysis for the aggregate data is contained in the “Aggregate” subdirectory.

- **AggregateBLSIPP.m:** This matlab program carries out the main calculations. The NIPA data on “intellectual property products” investment are from FRED; the download codes are reported in comments in the program.
- `mfp_tables_historical-2017-02-17.xls`: Contains the BLS data on private business sector TFP growth. The contribution from intellectual property products, which is netted out of TFP growth by the BLS, is added back in, in accordance with the model.

The idea output measure is TFP growth, by decade (and for 2000-2014 for the latest observation). For the years since 1950, this measure is the BLS Private Business Sector

multifactor productivity growth series, adding back in the contributions from R&D and IPP. For the 1930s and 1940s, we use the measure from Robert Gordon (2016). The idea input measure is gross domestic investment in intellectual property products from the National Income and Product Accounts, deflated by a measure of the nominal wage for high-skilled workers.

Figure 1 shows alternative measures of aggregate research effort, confirming the statement in the main text that our results are robust to how we measure aggregate research. In particular, the “NIPA IPP” series, which is the baseline series we report in the main text, and the “U.S.” measure of total researchers in full-time equivalents are very similar. The OECD and OECD+ series show that if we include broader measures of research effort, the decline in aggregate research productivity would be comparable in size or larger. These results are taken from the `AggregateBLS_SciEng.m` matlab program, and the underlying data are collected in `OECD-MSTI-TotalResearchers.xls`.

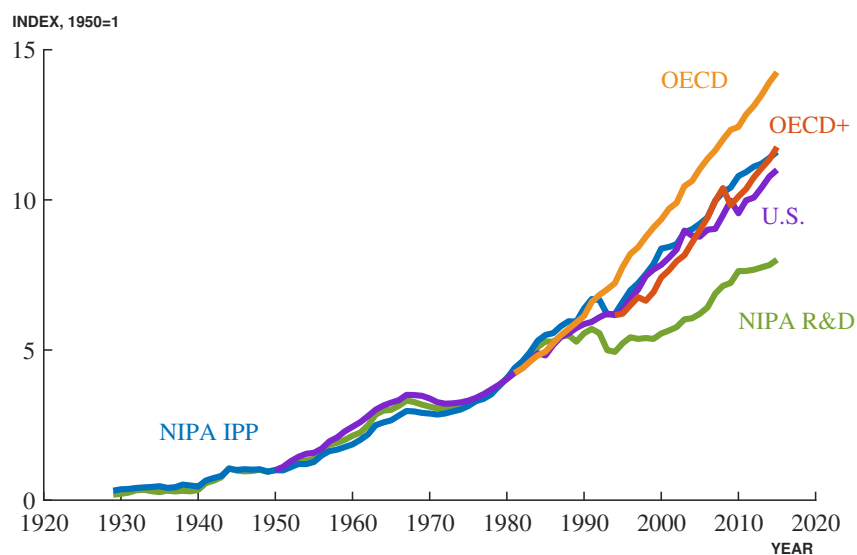
3. Moore’s Law

Our measurement of research spending related to Moore’s Law draws primarily on two sources. First, we use the Compustat database to obtain R&D spending for more than 35 multinational companies; we are grateful to Unni Pillai for his advice and preliminary data on semiconductor R&D.¹ Second, we use the PATSTAT database to obtain the fraction of each company’s patents that are in technology class “H01L” which is the class corresponding to semiconductors; we are grateful to Antoine Dechezlepretre for extensive help and computer code for extracting data from PATSTAT. Our various measures combine these data in different ways to create a measure of R&D relevant for Moore’s Law.

The spreadsheet “MooresLawRND-2018-01-08.xls” provides the basic background behind these calculations. The sheet labeled “Compustat” collates the Compustat R&D spending numbers with the (smoothed) patent shares. The sheet “PatentNarrow” provides our “narrow” measures — in which all firms research spending is weighted by their share of patents in the semiconductor class — while the sheet “PatentBroad” provides our “broad” measures — in which the research spending by focused companies

¹These data are supplemented in a few cases — for example for Siemens (thanks to Dietmar Harhoff) and Samsung (thanks to Jihee Kim) — by data from company annual reports.

Figure 1: Alternative Measures of Aggregate Research Effort



Note: The figure shows alternative measures of aggregate research effort. The “NIPA IPP” series is the main one reported in the paper; the “NIPA R&D” series includes only U.S. R&D expenditures, also as measured by the NIPA. Both of these series are deflated by the high-skilled wage series, as described in the main text. The other three series show measures of “Total Researchers (FTE)” from the OECD Main Science and Technology Indicators, <http://stats.oecd.org/ViewHTML.aspx?QueryId=58469#>. The U.S. line reports researchers in the United States; data before 1981 are taken from Jones (2002). The “OECD” line plots total researchers in OECD countries since 1981, showing a 3.4-fold increase since that year. The “OECD+” line adds researchers from China and Russia to the OECD measure and reveals a 1.9-fold increase between 1994 and 2015. For visual clarity, the OECD and OECD+ lines are normalized to the U.S. value in their starting years.

like Intel or Fairchild is all included, while the research spending by conglomerates like AT&T or IBM or Toshiba is weighted according to their semiconductor patent shares.

TFP growth in the “semiconductor and related device manufacturing” industry (NAICS 334413) is taken from the NBER/CES Manufacturing Industry Database, variable “dtfp5”; see Bartelsman and Gray (1996). We smooth TFP growth using an HP filter with smoothing parameter 400 and lag R&D by 5 years in computing research productivity. In addition to the narrow/broad split, we also alternately include and exclude R&D from semiconductor equipment manufacturers: equipment is captured in a separate 6-digit industry — and therefore is perhaps most naturally excluded from our analysis. Alternatively, the pricing of the semiconductor equipment may not fully capture the benefits of that equipment, in which case the R&D from semiconductor equipment manufacturing spills over into TFP growth in the semiconductor manufacturing sector.

The following matlab programs are used:

- **MooreLaw/IntelGraph.m:** This program produces the main results for Moore’s Law reported in the paper.
- **MooreLaw/SemiconductorTFP.m:** This program produces the main results for TFP growth in the semiconductor industry.
- **Patents/ReadPatentData.m:** This program reads the PATSTAT patent data and constructs the smoothed share of each firm’s patenting that is in the semiconductor class. These numbers show up in the Compustat tab of the main MooreLawRND*.xls spreadsheet.

4. Agricultural Innovations

The key files are contained in the “Seeds” subdirectory:

- **Seed data v6.xlsx:** This file contains the details of the data on seed yields and research spending.
- **SeedYields.m:** This matlab program carries out the main calculations that are reported in the paper.
- **AgIdeaPF.m:** TFP growth and research productivity for the agriculture sector as a whole. Both TFP growth and U.S. R&D spending for the agriculture sector as

a whole are taken from the U.S. Department of Agriculture's Economic Research Service. The TFP series is smoothed with an HP filter. Global R&D spending for agriculture is taken from Fuglie, Heisey, King, Day-Rubenstein, Schimmelpfennig, Wang, Pray and Karmarkar-Deshmukh (2011), Beintema, Stads, Fuglie and Heisey (2012), and Pardey, Chan-Kang, Beddow and Dehmer (2016). Nominal R&D spending is deflated by the wage for college graduates, as described earlier. The data are collected in the spreadsheet files "AgTFP v1.xlsx", USDA-ERS-ag_all_research.xls, and GlobalRND-Agriculture.xls.

We calculate idea input and output measures for agricultural crop yields in the United States for each of four crops: corn, soybeans, cotton, and wheat.

4.1. Idea output measure: crop yields

For each crop, we use realized average yields in the United States, measured in bushels or pounds harvested per acre planted. These data are provided by the U.S. Department of Agriculture (USDA) National Agricultural Statistics Service (NASS) (U.S. Department of Agriculture National Agricultural Statistics Service (2016)). We use yield figures provided annually at the national level. The URL to access the data is:

<https://quickstats.nass.usda.gov/#C57CA751-B131-3065-9F7C-E7DE08D92F87>

To obtain our final measure, we compute smoothed yields using an HP filter with a smoothing parameter of 400, then take an annualized average 5-year growth rate.

4.2. Idea input measure: seed R&D

For each crop, we calculate annual R&D expenditure in the United States directed at improving that crop's yields. Our data sources have three relevant dimensions: crop (corn, soybeans, cotton, and wheat), sector (public and private), and research area (biological efficiency, and crop protection and maintenance). For the private sector, we have measures of expenditure by research area that come aggregated over crops, so combine these with data on the share of a given research area devoted to a particular crop to produce an annual series of research area spending by crop. For the public sector, we have measures of total R&D by crop that come aggregated over research areas, so combine these with data on the share of a given crop's total R&D devoted to a

particular research area to produce an annual series of research area spending by crop. We sum across the two sectors to get an estimate for each crop-research area-year cell.

For the private sector, our measures of expenditure by research area were provided by Keith Fuglie of the U.S. Department of Agriculture's Economic Research Service. These series are an updated version of the series in Fuglie et al. (2011), and are annual from 1960 to 2015, in nominal dollars. The distribution of expenditure for seed efficiency by crop are taken from Perrin et al. (1983), Fernandez-Cornejo et al. (2004), Traxler et al. (2005), and Huffman and Evenson (2006). These provide shares for the years 1960, 1965, 1975, 1979, 1982, 1989, 1994, and 2001. No direct data are available on distribution of expenditure for crop protection by crop. As such, we took 3 related measures (crop protection sales shares from University of York (2016), the public crop protection shares described below, and the private seed efficiency shares described above) and took the average. We use linear interpolation where required to fill in missing years.

For the public sector, we begin with two comparable raw series of R&D by crop, covering different time periods, that each come aggregated over research areas. One is taken from annual versions of Table C from the U.S. Department of Agriculture National Institute of Food and Agriculture Current Research Information System (CRIS) Funding Summaries, and covers the years 1993-2015. The other is from Huffman and Evenson (2006), covering 1969, 1984, and 1997. The year 1997 is thus an overlapping year. CRIS figures are in nominal dollars; H&E figures are deflated by a price index. When we use this index to un-deflate the H&E series to get back nominal figures, the CRIS amounts for the overlapping year (1997) are close to 60% of the un-deflated H&E figures for all four crops. As such, we use the CRIS numbers for all the years available (1993-2015), and multiply the un-deflated H&E figures for the years 1969 and 1984 by this 'splicing factor' to get a consistent nominal series. The distribution of expenditure by research area for each crop is taken from Huffman and Evenson (2006). This provides shares for the years 1969, 1984, and 1997. We use linear interpolation where required to fill in missing years. We use the output from this methodology for spending on biological efficiency; for spending on biological efficiency and crop protection and maintenance combined, we use a new series provided by Huffman of absolute productivity-directed public research by crop for the years 1960-2009. This series is, as expected, very close

to the equivalent series generated using the methodology just outlined.

To obtain our final measure of idea inputs, we deflate the summed private and public annual series using a measure of the average annual earnings for people with 4 or more years of college, for reasons explained in section 3.

5. Medical Innovations

Programs and data for the disease measures are in the “Mortality” subdirectory, while those for the new molecular entities are in the “Pharma” subdirectory.

- **Cancer.m**, **BreastCancer.m**, **HeartDisease.m**: These are the main programs that carry out the calculations for the three diseases.
- **mortality.m**: This function is called to do the heavy lifting.
- **LifeExpectancy.m**: Create the life expectancy graph in Figure 7.
- **BasicLifeTable.m**: Reads the basic life tables from Mortality.org for all people.
- **BasicLifeTableWomen.m**: Reads the basic life tables from Mortality.org for women.
- **NMEGraph.m**: The basic program for generating the results for new molecular entities.
- **NME-Since1938.xls**: Data on new molecular entities since 1938, from <http://www.fda.gov/AboutFDA/WhatWeDo/History/ProductRegulation/SummaryofNDAApprovals> Also the R&D data from various issues of “Pharmaceutical Industry Profile”; see <http://www.phrma.org/sites/default/files/pdf/PhRMA%20Profile%202013.pdf>

Our measures of life expectancy and mortality from all sources by age come from the Human Mortality Database at <http://mortality.org>. To measure the percentage declines in mortality rates from cancer, we use the age-adjusted mortality rates for people ages 50 and over computed from 5-year survival rates, taken from the National Cancer Institutes Surveillance, Epidemiology, and End Results program at <http://seer.cancer.gov/>. For heart disease, we report the crude death rate in each year for people aged 55–64.

For our research input, we measure the number of scientific publications in PUBMED that have “Neoplasms” or “Breast Cancers” or “Heart Diseases”, as a MESH (Medical

Subject Heading) term. MESH is the National Library of Medicine's controlled vocabulary thesaurus. For more information on MESH, see <https://www.nlm.nih.gov/mesh/>. Our queries of the PUBMED data use the webtool created by the Institute for Biostatistics and Medical Informatics (IBMI) Medical Faculty, University of Ljubljana, Slovenia available at <http://webtools.mf.uni-lj.si/>.

6. Compustat Firm-Level Results

These programs and files are in the "Compustat" subdirectory.

- **Compustat-WRDS-2016-06-13.xlsx, Compustat-WRDS-2016-06-13.csv:** Basic data file downloaded from Compustat via WRDS.
- **CompustatRead.m:** Reads the downloaded data from Compustat-WRDS-2016-06-13.csv.
- **MasterCompustat.m:** The master program for the Compustat results, including robustness.
- **CompustatIdeaPF.m:** The basic program that does the heavy lifting, given a set of parameters and assumptions.
- **SetParameters.m:** Sets the baseline parameter values.
- **ShowParameters.m:** Reports the parameter values.
- **GDPDeflator.m:** Loads and saves the basic GDP Deflator used to deflate sales revenue and market cap.
- **compugrowthrate.m:** A function for computing various growth rates.

As a measure of the output of the idea production function, we use decadal averages of annual growth in sales revenue, market capitalization, employment, and revenue labor productivity within each firm. Sales revenue and market cap are deflated by the GDP implicit price deflator. We take the decade as our period of observation to smooth out fluctuations.

To measure the research input, we use a firm's spending on research and development from Compustat. This means we are restricted to publicly-listed firms that

report formal R&D, and such firms are well-known to be a select sample (e.g. disproportionately in manufacturing and large). We look at firms since 1980 that report non-zero R&D, and this restricts us to an initial sample of 15,128 firms. Our additional requirements for sample selection in our baseline sample are

1. We observe at least 3 annual growth observations for the firm in a given decade. These growth rates are averaged to form the idea output growth measure for that firm in that decade.
2. We only consider decades in which our idea output growth measure for the firm is positive (negative growth is clearly not the result of the firm innovating, and our framework cannot make sense of negative research productivity).
3. We require the firm to be observed (for both the output growth measure and the research input measure) for two consecutive decades. Our decades are the 1980s, the 1990s, the 2000s (which refers to the 2000-2007 period), and the 2010s (which refers to the 2010-2015 period); we drop the years 2008 and 2009 because of the financial crisis.

We relax many of these conditions in our robustness checks.

7. The Wage Series for Deflating R&D Spending

The program and data for this series are stored in the subdirectory “WageSci.”

- **WageEducation.m:** Reads the wage data and creates WageScientistData.mat, which is used in many other programs.

As a measure of the nominal wage in our empirical applications, we use mean personal income from the Current Population Survey for males with a Bachelor’s degree or more of education. These data are from Census Tables P18 and P19, available at <http://www.census.gov/topics/income-poverty/income/data/tables.html>. Prior to 1991, we use the series for “4 or more years of college.” For years between 1939 and 1967, we use the series Bc845 from the Historical Statistics for the U.S. Economy, Millennial Edition. Finally, for the aggregate research productivity calculations, we require a deflator from the 1930s. We extrapolate the college earnings series backward into the 1930s

using nominal GDP per person for this purpose. As an alternative, we have redone our results using nominal GDP per person as the deflator; this yielded broadly similar results.

References

- Bartelsman, Eric J. and Wayne Gray, "The NBER Manufacturing Productivity Database," Working Paper 205, National Bureau of Economic Research October 1996.
- Beintema, Nienke, Gert-Jan Stads, Keith Fuglie, and Paul Heisey, "ASTI Global Assessment of Agricultural R&D Spending," Technical Report, International Food Policy Research Institute 2012.
- Fernandez-Cornejo, Jorge et al., "The seed industry in US agriculture: An exploration of data and information on crop seed markets, regulation, industry structure, and research and development," Technical Report, United States Department of Agriculture, Economic Research Service 2004.
- Fuglie, Keith, Paul Heisey, John L King, Kelly Day-Rubenstein, David Schimmelpfennig, Sun Ling Wang, Carl E Pray, and Rupa Karmarkar-Deshmukh, "Research investments and market structure in the food processing, agricultural input, and biofuel industries worldwide," *USDA-ERS Economic Research Report*, 2011, (130).
- Gordon, Robert J., *The Rise and Fall of American Growth: The US Standard of Living since the Civil War*, Princeton University Press, 2016.
- Huffman, Wallace E and Robert E Evenson, *Science for agriculture: A long-term perspective*, John Wiley & Sons, 2006.
- Jones, Charles I., "Sources of U.S. Economic Growth in a World of Ideas," *American Economic Review*, March 2002, 92 (1), 220–239.
- Pardey, Philip G., Connie Chan-Kang, Jason M. Beddow, and Steven P. Dehmer, "Shifting Ground: Food and Agricultural R&D Spending Worldwide, 1960-2011," Technical Report, International Science and Technology Practice and Policy (InSTePP) Center, University of Minnesota 2016.
- Perrin, Richard K, KA Kunnings et al., "Some effects of the US Plant Variety Protection Act of 1970.," *North Carolina State University. Dept. of Economics and Business. Economics research report (USA). no. 46.*, 1983.

Traxler, Greg, Albert KA Acquaye, Kenneth Frey, and Ann Marie Thro, "Public sector plant breeding resources in the US: Study results for the year 2001," *USDA Cooperative State Research, Education and Extension Service*, 2005, pp. 1–7.

University of York, "The Essential Chemical Industry," 2016. [Online at <http://www.essentialchemicalindustry.org/>; accessed 8-September-2016].

U.S. Department of Agriculture National Agricultural Statistics Service, 2016. [Online at <https://quickstats.nass.usda.gov/>; accessed 8-September-2016].