

Big Data: Yet Another Buzzword or Actual Big Deal?

Robert J. Robbins

rjr8222@gmail.com

Biomedical Research Institutions Information Technology Exchange



© 2013, BRIITE

Biomedical Research Institutions Information Technology Exchange

11-13 December 2013

What Is **BIG DATA** (answer from the interwebs)



Big Data is a paranoid electronic music project from the Internet. The band was formed out of a mutual distrust for technology and The Cloud, despite a growing dependence on them. Their music explores the ever-increasing relationship between man and machine, and more specifically how the internet has reshaped the human experience. || **Alan Wilkis** [producer] +**Daniel Armbruster** [vocals]



Compute Information. Science, 332(6025), 60 -65. http://www.martinhilbert.net/WorldInfoCapacity.html



Compute Information. Science, 332(6025), 60 -65. http://www.martinhilbert.net/WorldInfoCapacity.html

Original Definition

F	plication Delivery Strategies META Group
at	File: 9 File: 9 Author: Doug Lar
	3D Data Management: Controlling Data Volume, Velocity, and Variety. Current business conditions and mediums are pushing traditional data management principles to their limits, giving rise to novel, more formalized approaches.
	META Trend: During 2001/02, leading enterprises will increasingly use a centralized data warehouse to det a common business vocabulary that improves internal and external collaboration. Through 2003/04, or swellty, and integration woes will be tempered by data profiling technologies (for generating metad

The effect of the e-commerce surge, a rise in merger/acquisition activity, increased collaboration, and the drive for harnessing information as a competitive catalyst is driving enterprises to higher levels of consciousness about how data is meaned at its nost basic tional, analytical, and collaborative consistencies, changing economic conditions have made this job more difficult. E-commerce, in particular, has exploded data management challenges along three dimensions: volumes, wele ty, and variety. In 2001/00

Original Definition

Volume: More data than extant systems can handle.

- Velocity: Data arriving at a rapid and accelerating pace.
 - Variety: Data in a multitude of formats, from a multitude of sources, all (somewhat) relevant to the question at hand.

Buzzword

Proof: Google

big data" - Google Search - Mozilla Firefox File Edit View Higtory Bookmarks Icols Help Thin data" - Google Search							
♦ ♦ ♥ ▲ https://www.google.com/search:	q="big data"&ie=utf-8&oe=utf-8∾	q=t&rls=org.mozilla:en-US:of1	ficial&client=firefox-a&channel=	np&source=hp	☆ ⊽ C	₹ Google	۶ 🖡
Google	"big data	a" 🗲					
	Web	Images	Maps	Shopping	News	More -	Search to
	About 11,1	00,000 res	sults (0.28	seconds)			
	What is I Ad www.	Big Data .sas.com/	<u>i? - Exter</u> Big-Data	nsive Insight ▼	s On <mark>Big</mark>	Data - SA	<u>S.com</u>
	SAS Softw	are has 1,	514 followe	ers on Google+			
	Big Data E	xplained -	SAS, Big [Data & Hadoop	- Try Visua	I Analytics D	emo
	Big Data	a in 2013 tableauso	- tablea	usoftware.co m/ big-data ▼	<u>om</u>		
	7 Things y	ou Need to	Do About	Big Data in 20	13. Get the	Free Article	
۲	Tahlaau Sr	offware ha	e 1 512 foll	owers on Good	10+		•

Proof: Google

"big data" - Google Search - Mozilla Firefox File Edit Yiew Higtory Bookmarks Tools Help Toni data" - Google Search + + + +							
♦ ♦ ♥ ▲ https://www.google.com/search:	="big data"8ie=utf-88ae=utf-88ae=t&rls=org.mozilla:en-US:official8client=firefox-a&channel=np8source=hp	۹ ا					
Google	"big data"						
	Web Images Maps Shopping News More -	Search to					
(About 11,100,000 results (0.28 seconds)						
	What is Big Data ? - Extensive Insights On Big Data - SAS.com Ad www.sas.com/ Big-Data ▼ Get the Free White Paper.						
	SAS Software has 1,514 followers on Google+						
	Big Data Explained - SAS, Big Data & Hadoop - Try Visual Analytics I)emo					
	Big Data in 2013 - tableausoftware.com						
	Ad www.tableausoftware.com/big-data -						
	7 Things you Need to Do About Big Data in 2013. Get the Free Article	ļ					
.	Tableau Software has 1 512 followers on Google+						

Proof: Amazon

🗧 🔶 📽 🛅 🖉 www.amazon.com/s/	ref=nb_sb_noss?url=search-alias%3Dstripbooks&fi	eld-keywords="big data"	☆ ⊽ C) <mark>।</mark>	▼ Google		۹ 🖡
amazon Prime Robert	's Amazon.com Today's Deals	Gift Cards Sell Help	>See the deals	Presented by Amaz	EALS	WEEK wards Visa Card
Shop by Search	Books 👻 "big data"		Go Hello, Robert Your Account -	Your 0 Prime -	Cart -	Wish List ⊸
Books Advanced Search New	w Releases Best Sellers The N	lew York Times® Best Sellers Children's	Books Textbooks Sell Your Books Bes	t Books of the Month	Deals in	n Books
Departments < Any Category	Books⇒ <mark>""big data""</mark>				Share	🗹 f 🎔
Books Computers & Technology (564)	Showing 1 - 12 of 1,130 Result	is 📃 Detail 🔛 Image		Sort by	Relevance	•
Data Mining (135)	Book Format					
Databases (240) Information Management (95) Data Warehousing (52)	Paperback	Hardcover Kindle E	Edition HTML			
Software Business (46)						
Modeling & Simulation (45) Information Theory (45) Education & Reference (233)		and Mayer-Schonberger, Vikto	at Will Transform How We Live, Wo or (Mar 14, 2013)	rk, and Think b	y Cukier,	Kenneth
Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38)	BIC	Big Data: A Revolution The and Mayer-Schonberger, Vikto	at Will Transform How We Live, Wo or (Mar 14, 2013) Price	rk, and Think by New	y Cukier, Used	Collectible
Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112)	BIC	Big Data: A Revolution That and Mayer-Schonberger, Vikto ★★★★★★ ♥ (183) Formats Kindle Edition Whispersync for Voice-ready	at Will Transform How We Live, Wo or (Mar 14, 2013) Price \$12.99	rk, and Think by New	y Cukier, Used	Kenneth Collectible
Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more		Big Data: A Revolution That and Mayer-Schonberger, Vikto ★★★★★★ ♥ (183) Formats Kindle Edition Whispersync for Voice-ready Hardcover Usually ships in 1 to 4 weeks	at Will Transform How We Live, Wo or (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53	New \$14.98	y Cukier, Used \$13.95	Kenneth Collectible \$36.95
Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more Amazon Prime \checkmark \checkmark <i>Prime</i> Eligible		Big Data: A Revolution The and Mayer-Schonberger, Vikto	at Will Transform How We Live, Wor or (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53	New \$14.98	y Cukier, Used \$13.95	Kenneth Collectible \$36.95
Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more Amazon Prime Mex Releases Last 30 days (69)	1. COCK INSIDE: BROCK DREVOLUTION INFERIMENTATIONA INFERIMENTATION INFERIMENTATION I	Big Data: A Revolution That and Mayer-Schonberger, Vikto →→→→→ ♥ (183) Formats Kindle Edition Whispersync for Voice-ready Hardcover Usually ships in 1 to 4 weeks Other Formats: Paperback; Audible Big Data For Dummies by H 2013)	at Will Transform How We Live, Wor or (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53 \$27.00 \$20.53	New New s14.98	y Cukier, Used \$13.95	Kenneth Collectible \$36.95 (Apr 2,
Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more Amazon Prime Amazon Prime Eligible New Releases Last 30 days (69) Last 90 days (226) Coming Soon (57)	1. COCK INSIDE: COCK INSIDE:	Big Data: A Revolution That and Mayer-Schonberger, Vikto ★★★★★ ♥ (183) Formats Kindle Edition Whispersync for Voice-ready Hardcover Usually ships in 1 to 4 weeks Other Formats: Paperback; Audible Big Data For Dummies by H 2013) ★★★★★ ♥ (14)	at Will Transform How We Live, Wor or (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53 e Audio Edition Hurwitz, Judith, Nugent, Alan, Halper, Fe	New \$14.98	y Cukier, Used \$13.95	Kenneth Collectible \$36.95 (Apr 2,

Proof: Amazon

Amazon.com; big data ; books	+					
🕨 🔶 📽 🔛 🖉 www.amazon.com/s	/ref=nb_sb_noss?url=search-alias%3Dstripbooks&	field-keywords="big data"	⊽ C 🖁 - Ga	ogle		۹ ا
amazon Prime Robert	t's Amazon.com Today's Deals	Gift Cards Sell Help	See the deals Pr	DAY DE	EALS	WEEK vards Visa Card
Shop by Search	h Books 👻 "big data"		Go Hello, Robert Your Your Account → Prin	ne - 🚺	Cart -	Wish List ⊸
Books Advanced Search Ne	w Releases Best Sellers The	New York Times® Best Sellers Children's	Books Textbooks Sell Your Books Best Bo	oks of the Month	Deals in	Books
Departments	Books > ""big data""				Share 칠	🛛 f 🏏
Books Computers & Technology (564)	Showing 1 - 12 of 1,130 Resu	Its 📃 Detail 🕞 Image		Sort by F	Relevance	•
Data Mining (135)	Dock Format					
Databases (240)						
Data Warehousing (52)	Paperback	Hardcover Kindle E	dition HTML			
e ()						
Software Business (46)						
Software Business (46) Modeling & Simulation (45) Information Theory (45)		Big Data: A Revolution Tha and Mayer-Schonberger, Vikto	t Will Transform How We Live, Work, a r (Mar 14, 2013)	and Think by	y Cukier, ł	Kenneth
Software Business (46) Modeling & Simulation (45) Information Theory (45) Education & Reference (233)	1. LOOK INSIDE!	Big Data: A Revolution Tha and Mayer-Schonberger, Vikto	It Will Transform How We Live, Work, r (Mar 14, 2013)	and Think by	y Cukier, ł	Kenneth
Software Business (46) Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38)	1. LOOK INSIDE	Big Data: A Revolution Tha and Mayer-Schonberger, Vikto	It Will Transform How We Live, Work, r (Mar 14, 2013) Price	and Think by New	y Cukier, ł Used	Collectible
Software Business (46) Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112)		Big Data: A Revolution Tha and Mayer-Schonberger, Vikto Constant (183) Formats Kindle Edition Whispersync for Voice-ready	t Will Transform How We Live, Work, a r (Mar 14, 2013) Price \$12.99	and Think by New	y Cukier, ł Used	Collectible
Software Business (46) Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more		Big Data: A Revolution Tha and Mayer-Schonberger, Vikto A A A A Y (183) Formats Kindle Edition Whispersync for Voice-ready Hardcover Usually ships in 1 to 4 weeks	t Will Transform How We Live, Work, a r (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53 <i>Aprime</i>	New \$14.98	y Cukier, ł Used \$13.95	Collectible \$36.95
Software Business (46) Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more Amazon Prime Scipping Eligible		Big Data: A Revolution Tha and Mayer-Schonberger, Vikto Constant (183) Formats Kindle Edition Whispersync for Voice-ready Hardcover Usually ships in 1 to 4 weeks Other Formats: Paperback; Audible	t Will Transform How We Live, Work, a r (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53 <i>Arime</i> e Audio Edition	New \$14.98	y Cukier, ł Used \$13.95	Collectible \$36.95
Software Business (46) Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more Amazon Prime Set State Set Set Set Set Set Set Set Set Set S	1. LOOK INSIDE DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTION DEPENDENTI	Big Data: A Revolution Tha and Mayer-Schonberger, Vikto Constant (183) Formats Kindle Edition Whispersync for Voice-ready Hardcover Usually ships in 1 to 4 weeks Other Formats: Paperback; Audible Big Data For Dummies by H 2013)	t Will Transform How We Live, Work, a r (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53 Prime Audio Edition	New \$14.98 and Kaufman	y Cukier, ł Used \$13.95	Kenneth Collectible \$36.95 Apr 2,
Software Business (46) Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more Amazon Prime Or Prime Eligible New Releases Last 30 days (69) Last 90 days (226)	1. LOOK INSIDE BACCONTON INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVOLUTION INTERVO	Big Data: A Revolution Tha and Mayer-Schonberger, Vikto ★★★★★ ♥ (183) Formats Kindle Edition Whispersync for Voice-ready Hardcover Usually ships in 1 to 4 weeks Other Formats: Paperback; Audible Big Data For Dummies by H 2013) ★★★★★ ♥ (14)	t Will Transform How We Live, Work, a r (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53 Advine Audio Edition	New \$14.98	y Cukier, ł Used \$13.95 I, Marcia (J	Collectible \$36.95 Apr 2,
Software Business (46) Modeling & Simulation (45) Information Theory (45) Education & Reference (233) Database Design (38) Business & Investing (423) Science & Math (112) + See more Amazon Prime Amazon Prime Eligible New Releases Last 30 days (69) Last 90 days (226) Coming Soon (57)	1. LOOK INSIDE REVOLUTION REVOLUTION REVOLUTION REVOLUTION REVOLUTION REVOLUTION REVOLUTION REVOLUTION REVOLUTION REVOLUTION	Big Data: A Revolution Tha and Mayer-Schonberger, Vikto ★★★★★ ♥ (183) Formats Kindle Edition Whispersync for Voice-ready Hardcover Usually ships in 1 to 4 weeks Other Formats: Paperback; Audible Big Data For Dummies by H 2013) ★★★★★ ♥ (14) Formats	t Will Transform How We Live, Work, a r (Mar 14, 2013) Price \$12.99 \$27.00 \$20.53 Arime e Audio Edition Rurwitz, Judith, Nugent, Alan, Halper, Fern Price	New New \$14.98 and Kaufman New	y Cukier, ł Used \$13.95 , Marcia ("	Collectible \$36.95 Apr 2,

Proof: Everybody Needs to Know















Big Data Analytics in Biomedical Research

PLUS: Privacy and Biomedical Research: Building a Trust Infrastructure

Winter 2011/2012

INSIDE:

Tapping the Brain: Decoding fMRI Personalized Cancer Treatment Leveraging Social Media For Biomedical Research And more

























Ultimate Proof: Dilbert



Proof: Even LOLCats



Not Only Buzzword Status



Big Deal
Truly humungous analyses

Truly humungous analyses

Datafication of reality

Truly humungous analyses

Datafication of reality

When $n \rightarrow all$

Truly humungous analyses Datafication of reality When $n \rightarrow all$

Analytical issues

Truly humungous analyses Datafication of reality When $n \rightarrow all$ **Analytical issues Policy Issues**

New Insights

Datafication of Reality

















When the first scientific expedition reached the site, 20 years after the event, they found hundreds of square miles of residual devastation. Although some firsthand descriptions of the event were obtained, these were based on twentyyear-old memories of people who has seen the event from a great distance. No photographs exist of the event as it happened.







Pretty Exciting, huh?



,12; hat

ubeek

> ive, .he

est-

rical

with

gen-

they

ed in

ratic

tion.

MN),

reated

illion tart, the seemingly mundane, yet _______ to or erringersu attorney stems higher on the agenda of organizations such as the G20. ■

Eyes and ears

Two explosions last week demonstrated the importance of global monitoring.

n 15 February, the town of Chelyabinsk in the Russian Ural Mountains had an unexpected visitor. A meteor streaked high above the city, briefly blinded commuters and then shattered thousands of windows with a series of ear-splitting explosions. The event was recorded on mobile phones and car-dashboard cameras across the region, and YouTube soon filled with Hollywood-style disaster videos of the fireball, replete with some very colourful Russian commentary.

Local residents were not the only ones to record the blast. More than a dozen monitoring stations around the globe captured the ultralow-frequency infrasound signal of the meteorite as it broke up in the atmosphere. The stations are part of a much larger network of sensors

21 FEBRUARY 2013 | VOL 494 | NATURE | 281

ublishers Limited. All rights reserved

the second the second for the spliceosome, yet group II splicing and spliceosomal splicing involve the same two phosphoryl-transfer reactions and generate the same reaction intermediates and products. The crystal structure of a group II intron⁶ revealed an active site with two catalytic metal ions coordinated to phosphate groups in the domain V helix of the RNA that are spatially and functionally equivalent to the phosphates that Fica et al. show are coordinated between the U6 snRNA and the pre-mRNA. The parallels between these molecules even extend to the stereochemistry of the oxygen atoms that serve as ligands for the two metal ions. Although it has long been known that the hairpin-loop structures of U6 and domain V are functionally similar, this work demonstrates the chemical equivalence between these two systems in their metal-ion coordination.

There are more than 100 proteins in the spliceosome, so what role do they have if RNA is the catalyst of pre-mRNA splicing? Some proteins are retained throughout the splicing process, whereas others associate SOLAR SYSTEM

Russian skyfall

The recent entry of a 20-metre-wide celestial rock into Earth's atmosphere offered both a spectacular show and a source of invaluable data that advance our understanding of high-velocity impacts. SEE LETTERS P.235 & P.238

NATALIA ARTEMIEVA

wo papers in this issue focus on the scientific reconstruction of the asteroid-impact event that was observed on 15 February 2013 in Chelyabinsk, Russia — the largest impact event on Earth since the Tunguska blast of 1908. Borovička *et al.*¹ (page 235) extract substantial scientific information from several low-quality video recordings. Meanwhile, by analysing infrasound airwaves and the brightness of the light

flash generated by the impact, as well as the damage caused to Earth's surface, Brown *et al.*² (page 238) estimate the total energy of the event to have been 400–600 kilotons of trinitrotoluene (1 kt of TNT is equivalent to 4.185×10^{12} joules of energy), and say that such events occur more often than previously thought.

There are three fundamental physical processes involved as meteoroids — asteroids' smaller counterparts — enter the atmosphere. First, atmospheric drag decelerates meteoroids from their initial velocity of about

202 | NATURE | VOL 503 | 14 NOVEMBER 2013

© 2013 Macmillan Publishers Limited. All rights reserved

Echos in the Da

A major paper in *Nature* reported an analysis of the Chelyabinsk impactor.

The primary source data were fifteen YouTube videos...

LETTER

doi:10.1038/nature12671

The trajectory, structure and origin of the Chelyabinsk asteroidal impactor

Jiří Borovička¹, Pavel Spurný¹, Peter Brown^{2,3}, Paul Wiegert^{2,3}, Pavel Kalenda⁴, David Clark^{2,3} & Lukáš Shrbený¹

Earth is continuously colliding with fragments of asteroids and comets of various sizes. The largest encounter in historical times occurred over the Tunguska river in Siberia in 1908, producing^{1,2} an airburst of energy equivalent to 5 15 megatons of trinitrotoluene (1 kiloton of trinitrotolucne represents an energy of 4.185×10^{12} joules). Until recently, the next most energetic airburst events occurred over Indonesia³ in 2009 and near the Marshall Islands⁴ in 1994, both with energies of several tens of kilotons. Here we report an analysis of selected video records of the Chelyabinsk superbolide⁵ of 15 February 2013, with energy equivalent to 500 kilotons of trinitrotoluene, and details of its atmospheric passage. We found that its orbit was similar to the orbit of the two-kilometre-diameter asteroid 86039 (1999 NC43), to a degree of statistical significance sufficient to suggest that the two were once part of the same object. The bulk strength-the ability to resist breakage-of the Chelyabinsk asteroid, of about one megapascal, was similar to that of smaller meteoroids6 and corresponds to a heavily fractured single stone. The asteroid broke into small pieces between the altitudes of 45 and 30 kilometres, preventing more-serious damage on the ground. The total mass of surviving fragments larger than 100 grams was lower than expected7. The data for Tunguska are limited to tree damage and records of

esimic or a coustic wave at large distances. The Indonesia and Marshall Islands impacts were detected only by distant infrassoic station

aßG.

The Chelyabinsk impact occurred unexpectedly over a relatively densely populated Russian region during sunrise on 15 February 2013. The superbolide (an extremely bright meteor) generated a damaging air blast wave. An 8-m-wide hole in the ice of Lake Chebarkul, 70 km west of Chelvabinsk, was reported shortly after the event. Thousands, of small meteorites of total mass >100 kg, classified as LL5 ordinary chondrites, were found in the areas south-southwest or Che Here we determine the bolide trajectory and orbit and describe the ablation process of the asteroid. The main data for these analyses were 15 bolide videos publicly available on the internet (Extended Data Table 1). 🛀 calibrated these videos with wide-field stellar imagery. Details, procedure, which was based on the least squares memod,", are given in Supplementary Information. The trajectory and speed of the bolide are presented in Table 1. The observed low deceleration provides an extreme lower limit of 10° kg for the mass of the body. The measured energy⁵ and speed provide a best estimate of the mass of $\sim 1.2 \times 10^7$ kg, correspond-

ing to a diameter of ~19 m assuming a bulk density of $3,300 \, \text{kg}\,\text{m}^{-3}$. The pre-impact orbit (Table 2) is consistent with an origin in the main astroid belt, most probably in the inner main belt near the v_s socular resonance. We integrated the orbit and 1,000 test particles within the orbital uncertaintics (a probability cloud) 2,000 years into the past. The astroid spent the six weeks he fore impact within an elongation of 45°

from the Sun, a region of the sky inaccessible to ground-based telescopes. At earlier times, the asteroid was always too faint to be seen when some portion of the probability cloud was in the field of view of existing asteroid surveys. We note that the 2.2-km-diameter¹⁵ near-Earth asteroid 86039 (1999 NC43) of spectral type Q16 (corresponding to ordinary chondrites) has a very similar orbit, with very low dissimilarity criteria, D = 0.050 (ref. 17) and D' = 0.018 (ref. 18), relative to Chelyabinsk asteroid. Though this does not provide an unequivocal dynamical link, such a close match is unlikely statistically. We expect 227 near-Earth asteroids brighter than 86039 to exist¹⁹. Selecting at random from the expected distribution of near-Earth asteroids²⁰, it takes an average of 6×10^5 draws before selecting one with a smaller D value, and more than 3×10^6 draws before selecting one with a smaller D' value. Because 227/600,000 and 227/3,000,000 are equivalent to 1:2,600 and 1:13,000, respectively, we conclude that there is an approximately 1:10,000 chance that the proximity of these orbits is due purely to chance. The two orbits have maintained two intersection points over the past 2,000 years, one near perihelion and one near aphelion (Extended Data Fig. 1). The minimum velocity kick required to eject the Chelvabinsk. asteroid from 86039 is 0.7 km s⁻¹ (aphelion) or 2 km s⁻¹ (perihelion). This ejection velocity is consistent with a collision with another asteroid (which would provide a kick of a few kilometres per second). The frammation wring atmospheric entry was studied using the

Table 1 | Trajectory of the Chelyabinsk superbolide

Time (s)	Longitude (°)	Latitude (")	Height (km)	Speed (kms ⁻¹)
1.07	64.477	54.454	95.0	19.03
6.97	G2.888	54.664	60.0	19.05
10.46	61.933	54.780	40.0	19.03
12.24	61.442	54.837	30.0	18.9
13.18	61.193	54.864	25.0	18.0
14.18	60.943	54.892	20.0	14.2
5.17	60.802	54.907	17.2	6
		Fragment F1		
14.32	60.945	54,893	20.0	13.5
16.04	60.704	54.922	15.0	6.4
17.80	60.5883	54.9361	12.57	3.2

Time are occrease outside comparingly 2020/2011. Coordinate are given in the WCSB4 evol disternative concentration of the Cart's standard. At the bagining the lacose each out but profile was user than excellentian in the antimophenic drug. The heighting specific(14):10.13km s⁻¹) environed and the transformation of the transfo

¹Astronomical Institute Academy of Sciences of the Dresh Republic, 07-051 650 md Rejus Drech Republic. ¹Department of Physics and Ast mcomy. University of Western Dictario, London, Dictario, NGA XX7, Canada. ²Centre for Kinetary Science and Exploration. University of Western Ontario, London, Unitario, Kada Sciences of the Cascin Republic, Y Holeson Visio, et al., 42:100 Paral & Cascin Republic.

, and c.

14 NOVEMBER 2013 | VOL 503 | NATURE | 235 (2013 Macmilian Publishers Limited. All rights received

Echos in the Dat

A major paper in *Science* reported an analysis of the Chelyabinsk impactor.

Some critical data were taken from YouTube videos...

Chelyabinsk Airburst, Damage Assessment, Meteorite Recovery, and Characterization

Olga P. Popova,¹ Peter Jenniskens,^{2,3}, Vachestav Emet'yanenko,⁴ Anna Kartashova,⁴ Eugeny Biryukov,⁵ Sergey Khaibrakhmanov,⁶ Valery Shuvalov,¹ Yurij Rybnov,¹ Alexandr Dudorov,⁶ Victor I. Grokhovsky,⁷ Dmitry D. Badyukov,⁹ Qing-Zhu Yin,⁹ Peter S. Gural,² Jim Albers,² Mikael Granvik,¹⁰ Läslo G. Evers,^{11,12} Jacob Kuiper,¹¹ Vladimir Kharlamov,¹ Andrey Solovyov, 13 Yuri S. Rusakov, 14 Stanislav Korotkiy, 15 Ilya Serdyuk, 16 Alexander V. Korochantsev,⁸ Michail Yu. Larionov,⁷ Dmitry Glazachev,¹ Alexander E. Mayer,⁶ Galen Gisler,¹⁷ Sergei V. Gladkovsky,¹⁸ Josh Wimpenny,⁹ Matthew E. Sanborn, Akane Yamakawa,⁹ Kenneth L. Verosub,⁹ Douglas J. Rowland,¹⁹ Sarah Roeske,¹ Nicholas W. Botto,⁹ Jon M. Friedrich,^{20,21} Michael E. Zolensky,²² Loan Le,^{23,22} Daniel Ross,^{23,22} Karen Ziegler,²⁴ Tomoki Nakamura,²⁵ Insu Ahn,²⁵ Jong Ik Lee,²⁶ Qin Zhou,^{27,28} Xian-Hua Li,²⁸ Qiu-Li Li,²⁸ Yu Liu,²⁸ Guo-Qiang Tang,²⁸ Takahiro Hiroi,²⁹ Derek Sears,⁵ Ilya A. Weinstein,² Alexander S. Vokhmintsey," Alexei V. Ishchenko," Phillipe Schmitt-Kopplin, 50,5 Norbert Hertkorn,³⁰ Keisuke Nagao,³² Makiko K. Haba,³² Mutsumi Komatsu,³³ Takashi Mikouchi, 34 (the Chelyabinsk Airburst Consortium)

The asteroid impact near the Russian city of Chelyabinsk on 15 February 2013 was the largest airburst on Earth since the 1908 Tunguska event, causing a natural disaster in an area with a population exceeding one million. Because it occurred in an era with modern consumer electronics, field sensors, and laboratory techniques, unprecedented measurements were made of the impact event and the meteoroid that caused it. Here, we document the account of what happened, as understood now, using comprehensive data obtained from astronomy, planetary science, geophysics, meteorology, meteoritics, and cosmochemistry and from social science surveys. A good understanding of the Chelyabinsk incident provides an opportunity to calibrate the event, with implications for the study of near-Earth objects and developing hazard mitigation strategies for planetary protection.

energy (km⁻¹).

- helyabinsk Ublast - inner - certica - 1) The C-t-



Fig. 1. Meteoroid fragmentation stages in video taken by A. Ivanov in Kamensk-Uralskiy. (A) Fireball just before peak brightness, at the moment when camera gain was first adjusted. (B) End of main disruption. (C) Onset of secondary disruption. (D) End of secondary disruption; main debris cloud continues to move down. (E) Two main fragments remain. (F) Single fragment remains. (G) Thermally emitting debris doud at rest with atmosphere. (H) Final fragment continues to penetrate. Meteor moved hehind distant lamp posts. (Land 1) Detail of the thermal emission from a photograph by Mr. Dudarev (I) and M. Ahmetvaleev (J), after sky subtraction with high-pass filter and contrast enhancement. Altitude scale is uncertain by +0.7 km.

trajectory were pushed inwards, and suspended ceilings were sucked down above broken windows (fig. \$36G). There was no structural damage to buildings, other than a statue of Pushkin inside the local library, cracked by a blown-out window frame. Cracks in walls were documented in nearby Baturinsky and Kalachevo.

Electrophonic sounds were heard (SM section 1.6), but there was no evidence of an electromagnetic pulse (EMP) under the track in neighboring Emanzhelinka. Due to shock-wave-induced vibrations, electricity and cell phone connectivity was briefly halted in the Kunashaksky district at the far northern end of the damage area. The gas supply was briefly interrupted in some districts because of valves reacting to the vibrations.

People found it painful to look at the bright fireball, but glancing away prevented lasting eye damage. Of 1113 respondents to an Internet survey who were outside at the time, 25 were sun- The location of the meteorites is consistent with hume \$ 42.2%), 315 feb hot (28%) and 415 and an orthwest winds of 5 to 15 (1)\$ (fig. \$24)

cording to $r^{-3.2 \pm 0.5}$ (SM section 2.1). The majority of injuries (1210) took place in the densely populated Chelyabinsk city, but the highest fraction of people asking for assistance was near the trajectory track in the Korkinsky district (0.16%)

Meteorite Recovery

Shock radiation contributed to surface heating and ablation but did not completely evaporate all fragments of Chelvabinsk, unlike in the case of Tunguska (3). Meteorites of ~0.1 g fell near Aleksandrovka close to the point of peak brightness, masses of ~100 g fell further along the trajectory near Deputatskiy, and at least one of 3.4 kg fell near Timiryazevskiy. One hit the roof of a house in Deputatskiy (fig. \$46). Fallingsphere models suggest that they originated at 32to 26-km altitude (fig. S52), where the meteor model shows rapid fragmentation (fig. S18C).

RESEARCHARTIC

efficiently decelerated, avoiding the transfer of momentum to lower altitudes and resulting in less damage when the blast wave reached the ground.

Damage Assessment

In the weeks after the event, 50 villages were visited to verify the extent of glass damage. The resulting map (Fig. 3) demonstrates that the shock wave had a cylindrical component, extending furthest perpendicular to the trajectory. There was little coherence of the shock wave in the forward direction, where the disturbance was of long duration, shaking buildings and making people run outside, but causing no damage

¹Institute for Dynamics of Geospheres of the Russian Academ of Sciences, Leninsky Prospect 38, Building 1, Moscow, 119334, Russia, ²SETI Institute, 189 Bernardo Avenue, Mountain View CA 94043, USA, ³NASA Ames Research Center, Moffett Field Mail Stop 245-1. CA 94035. USA. "Institute of Astronomy of the Russian Academy of Sciences, Pyatnitskaya 48, Moscow, 119017, Russia. ⁵Department of Theoretical Mechanics, South Ural State University, Lenin Avenue 76, Chelyabinsk, 454080, Russia, ⁶Chelvabinsk State University, Bratvey Kashirinyh Street 129, Chelyabinsk, 454001, Russia. ⁷Institute of Physics and Tech nology, Ural Federal University, Mira Street 19, Yekaterinburg 620002 Rossia Vemarkky institute of Geochemistry and Analytical Chemistry of the RAS, Kosygina Street 19, Moscow, 119991, Russia. ⁹Department of Earth and Planetary Sciences, University of California at Davis, Davis, CA 95616, USA, ¹⁰Department of Physics, University of Helsinki, P.O. Box 64, 00014 Helsinki, Finland. 12 Koninklijk Nederlands Meteorologisch Instituut, P.O. Box 201, 3730 ÅF De Bilt, Netherlands. ¹² Department of Geoscience and Engineering, Faculty of Civil Engineering and Geosciences, Delft University of Technology, P.O. Box 5048, 2600 GA Delft, Netherlands. 13 Tomsk State University, Lenina Prospect 36, Tomsk, 634050, Russia. 14Research and Production Association "Typhoon," Floor 2, 7 Enock Street



58

New Insights

Analysis Issues When $n \rightarrow all$





Excerpts:

At its core, big data is about predictions.

Big data's ascendancy represents three shifts in the way we analyze information that transform how we understand and organize society.



Excerpts:

- 1. In this new world we can analyze far more data. In some cases we can even process all of it relating to a particular phenomenon. ... the need for sampling is an artifact of a period of information scarcity, a product of the natural constraints on interacting with information in an analog era.
- 2. Looking at vastly more data also permits us to loosen up our desire for exactitude. ... It is a trade-off: with less error from sampling we can accept more measurement error.



Excerpts:

3. These two shifts lead to a third change: a move away from the age-old search for causality. As humans we have been conditioned to look for causes, even though searching for causality is often difficult and may lead us down the wrong paths. In a big-data world, by contrast, we won't have to be fixated on causality; instead we can discover patterns and correlations in the data that offer us novel and invaluable insights. The correlations may not tell us precisely *why* something is happening, but they alert us *that* it is happening.



Excerpts:

Sampling is an outgrowth of an era of informationprocessing constraints, when people were measuring the world but lacked the tools to analyze what they collected.

Yet sampling comes with the cost of the has long been acknowledged but shunted aside. It loses detail. In some

cases there is no other way but to sample. In many areas however, a shift is taking place from collecting some data, to gathering as much as possible, and if feasible, getting everything:

N = all



N = all:

When it is possible to collect data on all of the objects of interest, profound changes in analytical methods are required.

In addition to changes in analytical methods, even more profound policy issues begin to arise.

New Insights

Policy Issues When $n \rightarrow all$



Policy Issues



Policy Issues

We all emit a constant stream of data exhaust that can be used to monitor our daily activities... The Economist

NOVEMBER 16TH-22ND 2013

Economist.com

China paves the way for reform Climate, storms and the Philippines The lure of foxy faces Europe's far-right alliance Buddhism and business

Every step you take

Google Glass, ubiquitous cameras and the threat to privacy

Finding a needle in haystack is fairly hard ...

Finding a needle in haystack is fairly hard ...

unless you have a really big electromagnet.

Connecting data to one individual on Earth is fairly hard ...

Connecting data to one individual on Earth is fairly hard ...

unless you have access to several big-data data sets, AND you know just a few, independent attributes of the person.






Where on Earth was this picture taken?

Method:

- Download JPG
- Extract date from EXIF information
- Google on date & waterspout | tornado

Result:

 Pictures were taken: Adriatic sea, Croatia, between islands Hvar and Brac just above village Murvica beach on island Brac.









True Story:

- I met a guy on a plane.
- Didn't get his name.
- Wished I had...

What I did get:

- He did IT consulting for pharmas
- He was German
- Seattle was his home airport
- **PROSIS (?) is his company name**
- His first name sounds like Hakim
- He knows Tim Jenkins

What I did when I got home:

- Opened Jenkins' Linked-In page
- Started looking at all contacts
- Found the name Achim Reeb
- Achim = Hakim, maybe ??
- Google "achim reeb"
- BINGO !

What I did when I got home:

- Opened Jenkins' Linked-In page
- Started looking at all contacts
- Found the name Achim Reeb
- Achim = Hakim, maybe ??
- Google "achim reeb"
- BINGO !



Hypothetical Story:

I met another guy on a plane ...









Hypothetical Story:

I met yet another guy on a plane ...





	arks <u>T</u> ools <u>H</u> elp	
8 nih informatics warren - Google Sear	rch +	
♦ ♦ % Attps://www	r. google.com/search?q=CIO&ie=utf-8&aq=t&rls=org.mozilla:en-US:official&client=firefox-a&channel=np&source=hp#channel=np&q=nih+inform 🏠 🤊 C 🕽 🗧 Google	₽ ♣ 🏫
		<u> </u>
Coorde	nih informatics warren	Q
Charles and		
	Web Imagaa Mana Channing Daaka Mara z Saarah taala	
	web images maps Shopping Books more Search tools	
	About $40,600,000$ require $(0.27,accorde)$	
	About 10,000,000 Tesuits (0.57 Seconds)	
	Dr Warren Kibbe to Head the NCI Center for Biomedical Informatics	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic ▼ Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic ▼ Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics , Warren has a the	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics , Warren has a the Director of the National Cancer Institute, National Institutes of Health .	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics, Warren has a the Director of the National Cancer Institute, National Institutes of Health.	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics , Warren has a the Director of the National Cancer Institute, National Institutes of Health . NCI BioMedical Informatics Blog NCIP Home - National Institute	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics , Warren has a the Director of the National Cancer Institute, National Institutes of Health . NCI BioMedical Informatics Blog NCIP Home - National Institute ncip.nci nih goy/blog/	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics , Warren has a the Director of the National Cancer Institute, National Institutes of Health . NCI BioMedical Informatics Blog NCIP Home - National Institute ncip.nci. nih .gov/blog/ Sep 27, 2013 - Lam pleased to appounce that Dr. Warren Kibbe has accepted my offer	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics , Warren has a the Director of the National Cancer Institute, National Institutes of Health . NCI BioMedical Informatics Blog NCIP Home - National Institute ncip.nci. nih .gov/blog/ Sep 27, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer to serve as the Director of the NCI Center for Biomedical Informatics and	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics , Warren has a the Director of the National Cancer Institute, National Institutes of Health . NCI BioMedical Informatics Blog NCIP Home - National Institute ncip.nci. nih .gov/blog/ Sep 27, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer to serve as the Director of the NCI Center for Biomedical Informatics and	
	ncip.nci.nih.gov//dr-warren-kibbe-to-head-the-nci-center-for-biomedic Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics, Warren has a the Director of the National Cancer Institute, National Institutes of Health. NCI BioMedical Informatics Blog NCIP Home - National Institute ncip.nci.nih.gov/blog/ Sep 27, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer to serve as the Director of the NCI Center for Biomedical Informatics and	
	ncip.nci.nih.gov//dr-warren-kibbe-to-head-the-nci-center-for-biomedic • Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics, Warren has a the Director of the National Cancer Institute, National Institutes of Health. NCI BioMedical Informatics Blog NCIP Home - National Institute ncip.nci.nih.gov/blog/ • Sep 27, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer to serve as the Director of the NCI Center for Biomedical Informatics and NCIP Home - National Institutes of Health	
	ncip.nci. nih .gov//dr- warren -kibbe-to-head-the-nci-center-for-biomedic • Jul 8, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer Thanks to his familiarity with NCI's past efforts in informatics , Warren has a the Director of the National Cancer Institute, National Institutes of Health . NCI BioMedical Informatics Blog NCIP Home - National Institute ncip.nci. nih .gov/blog/ • Sep 27, 2013 - I am pleased to announce that Dr. Warren Kibbe has accepted my offer to serve as the Director of the NCI Center for Biomedical Informatics and NCIP Home - National Institutes of Health ncip.nci. nih .gov/ •	



SCIENTISTS EXPOSE NEW VULNERABILITIES IN THE SECURITY OF PERSONAL GENETIC INFORMATION



Using only a computer, an Internet connection, and publicly accessible online resources, a team of Whitehead Institute researchers has been able to identify nearly 50 individuals who had submitted personal genetic material as participants in genomic studies.

GenomeHacking

Yaniv Erlich

Whitehead Institute for Biomedical Research

(yaniv@wi.mit.edu)



Yaniv Erlich Whitehead Institute for Biomedical Research

in histone pre-mRNA 3'-end processing (fig. S14).

. H. Durthan Et ung 3094 (2006).

.... IL SCI. U.J.A. 103,

winug. 2012; dec. All novem. 12 10.1126/science.1228705

Identifying Personal Genomes by **Surname Inference**

Melissa Gymrek,^{1,2,3,4} Amy L. McGuire,⁵ David Golan,⁶ Eran Halperin,^{7,8,9} Yaniv Erlich¹*

Sharing sequencing data sets without identifiers has become a common practice in genomics. Here, we report that surnames can be recovered from personal genomes by profiling short tandem repeats on the Y chromosome (Y-STRs) and querying recreational genetic genealogy databases. We show that a combination of a surname with other types of metadata, such as age and state, can be used to triangulate the identity of the target. A key feature of this technique is that it entirely relies on free, publicly accessible Internet resources. We quantitatively analyze the probability of identification for U.S. males. We further demonstrate the feasibility of this technique by tracing back with high probability the identities of multiple participants in public sequencing projects.

urnames are paternally inherited in most human societies, resulting in their co-Segregation with Y-chromosome haplotypes (1-5). Based on this observation, multiple genetic genealogy companies offer services to reunite distant patrilineal relatives by genotyping a few dozen

¹Whitehead Institute for Biomedical Research, 9 Cambridge Center, Cambridge, MA 02142, USA. ²Harvard–Nassachusetts Institute of Technology (MIT) Division of Health Sciences and Technology, MIT, Cambridge, MA 02139, USA. ³Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA. ⁴Department of Molecular Biology and Diabetes Unit, Massachusetts General Hospital, Boston, MA 02114, USA. ^SCenter for Medical Ethics and Health Policy, Baylor College of Medicine, Houston, TX 77030, USA. ⁶Department of Statistics and Operations Research, Tel Aviv University, Tel Aviv 69978, Israel. ⁵School of Computer Science, Tel Aviv University, Tel Aviv 69978, Israel. ⁸Department of Molecular Microbiology and Biotechnology, Tel-Aviv University, Tel Aviv 69978, Israel. 9The International Computer Science Institute, Berkeley, CA 94704, USA.

*To whom correspondence should be addressed. E-mail: yaniv@wi.mit.edu

highly polymorphic short tandem repeats across the Y chromosome (Y-STRs). The association between sumames and haplotypes can be confounded by nonpaternity events, mutations, and adoption of the same sumame by multiple founders (5). The genetic genealogy community addresses these barriers with massive databases that list the test results of Y-STR haplotypes along with their corresponding sumames. Currently, there are at least cight databases and numerous sumame project Web sites that collectively contain hundreds of thousands of surname-haplotype records (table S1).

The ability of genetic genealogy databases to breach anonymity has been demonstrated in the past. In a number of public cases, male adoptees and descendants of anonymous sperm donors used recreational genetic genealogy services to genotype their Y-chromosome haplotypes and to search the companies' databases (6-9). The genetic matches identified distant patrilineal relatives and pointed to the potential sumames of their biological fathers.

By combining other pieces of demographic information, such as date and place of birth, they fully exposed the identity of their biological fathers. Lunshof et al. (10) were the first to speculate that this technique could expose the full identity of participants in sequencing projects. Gitschier (11) empirically approached this hypothesis by testing 30 Y-STR haplotypes of CFU participants in these databases and reported that potential surnames can be detected. [CEU participants are multigenerational families of northern and western European ancestry in Utah who had originally had their samples collected by CEPH (Centre d'Etude du Polymorphisme Humain) and were later reconsented to participate in the HapMap project.] However, these surnames could match thousands of individuals, and the study did not pursue full re-identification at a single-person resolution.

Our goal was to quantitatively approach the question of how readily sumame inference might be possible in a more general population, apply this approach to personal genome data sets, and demonstrate end-to-end identification of individuals with only public information. We show that full identities of personal genomes can be exposed via surname inference from recreational genetic genealogy databases followed by Internet searches. In all cases in which individuals were studied who had donated DNA samples, the informed consent statements they had signed stated privacy breach as a potential risk and the data usage terms did not prevent re-identification. Representatives of relevant organizations that funded the original studies were notified and confirmed the compliance of this study with their guidelines (12).

As a primary resource for surname inference, we focused on Ysearch (www.ysearch.org) and

www.sciencemag.org SCIENCE VOL 339 18 JANUARY 2013

http://www.briite.org

321

3094 (2006).

IL. SCI. U.J.A. 100,

winug, 2012; dcc. All putero, 12 10.1126/science.1228705

in histone pre-mRNA 3'-end processing (fig. S14).

Identifying Personal Genomes by **Surname Inference**

Melissa Gymrek,^{1,2,3,4} Amy L. McGuire,⁵ David Golan,⁶ Eran Halperin,^{7,8,9} Yaniv Erlich¹*

without identifiers has become a common practice in genomics. sovered from personal genomes by profiling short tandem Sharing sequencing data se ving recreational genetic genealogy databases. Here, we report that surna her types of metadata, such as age and state, repeats on the Y chromosome (We show that a combination of a surname can be used to triangulate the identity of the arget. A key feature of this technique is that it entirely relies on free, publicly accessible Internet resources. We quantitatively analyze the probability of identification for U.S. males. We further demonstrate the feasibility of this technique by tracing back with high probability the identities of multiple participants in public sequencing projects.

urnames are paternally inherited in most human societies, resulting in their co-

segregation with Y-chromos (1-5). Based on this observation, genealogy companies offer service tant patrilineal relatives by genotyp

¹Whitehead Institute for Biomedical Re Center, Cambridge, MA 02142, USA. ²H Institute of Technology (MIT) Division of Technology, MIT, Cambridge, MA 02139, U.

ical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA. ⁴Department of Molecular Biology and Diabetes Unit, Massachusetts General Hospital, Boston, MA 02114, USA. ^SCenter for Medical Ethics and Health Policy, Baylor College of Medicine, Houston, TX 77030, USA. ⁶Department of Statistics and Operations Research, Tel Aviv University, Tel Aviv 69978, Israel. ⁵School of Computer Science, Tel Aviv University, Tel Aviv 69978, Israel. ⁸Department of Molecular Microbiology and Biotechnology, Tel-Aviv University, Tel Aviv 69978, Israel. 9The International Computer Science Institute, Berkeley, CA 94704, USA.

*To whom correspondence should be addressed. E-mail: yaniv@wi.mit.edu

highly polymorphic short tandem repeats across the Y chromosome (Y-STRs). The association bed hanlotunes can be confounded

By combining other pieces of demographic information, such as date and place of birth, they fully exposed the identity of their biological fathers. Lunshof et al. (10) were the first to speculate that this technique could expose the full identity of participants in sequencing projects. Gitschier (11) empirically approached this hypothesis by testing 30 Y-STR haplotypes of CFU participants in these databases and reported that potential surnames can be detected. [CEU participants are multigenerational families of northern and western European ancestry in Utah who had originally had their samples collected by CEPH (Centre d'Etude du Polymorphisme Humain) and were later reconsented to participate in the HapMap project.] However, these surnames could match thousands of individuals, and the study did not pursue full re-identification at a single-person resolution.

Our goal was to quantitatively approach the

A key feature of this technique is that it entirely relies on free, publicly accessible Internet resources

sands of surname-haplotype records (table S1). The ability of genetic genealogy databases to

breach anonymity has been demonstrated in the past. In a number of public cases, male adoptees and descendants of anonymous sperm donors used recreational genetic genealogy services to genotype their Y-chromosome haplotypes and to search the companies' databases (6-9). The genetic matches identified distant patrilineal relatives and pointed to the potential sumames of their biological fathers.

searches. In all cases in w studied who had donated DNA samples, the informed consent statements they had signed stated privacy breach as a potential risk and the data usage terms did not prevent re-identification. Representatives of relevant organizations that funded the original studies were notified and confirmed the compliance of this study with their guidelines (12). As a primary resource for surname inference,

al

et

321

we focused on Ysearch (www.ysearch.org) and

www.sciencemag.org SCIENCE VOL 339 18 JANUARY 2013

http://www.briite.org

And yet...

NIH GWAS policy, soon to be included in the more general NIH genomic-data sharing policy, requires that institutions stipulate that:

And yet...

NIH GWAS policy, soon to be included in the more general NIH genomic-data-sharing policy, requires that institutions stipulate that:

Data should be de-identified according to the following criteria: the identities of data subjects cannot be readily ascertained or otherwise associated with the data by the repository staff or secondary data users (45 C.F.R. 46.102(f)); the 18 identifiers enumerated at section 45 C.F.R. 164.514(b)(2) (the HIPAA Privacy Rule) are removed; and the submitting institution has no actual knowledge that the remaining information could be used alone or in combination with other information to identify the subject of the data.

And yet...

NIH GWAS policy, soon to be included in the more general NIH genomic-data-sharing policy, requires that institutions stipulate that:

Data should be de-identified according to the following criteria: the identities of data subjects cannot be readily ascertained or otherwise associated with the data by the repository staff or secondary data users (45 C.F.R. 46.102(f)); the 18 identifiers enumerated at section 45 C.F.R. 164.514(b)(2) (the HIPAA Privacy Rule) are removed; and the submitting institution has no actual knowledge that the remaining information could be used alone or in combination with other information to identify the subject of the data.

And yet...

NIH GWAS policy, soon to be included in the more general NIH genomic-data-sharing policy, requires that institutions stipulate that:

Data should be de-identified according to the following criteria: the identities of data subjects cannot be readily ascertained or otherwise associated with the data by the repository staff or secondary data users (45 C.F.R. 46.102(f)); the 18 identifiers enumerated at section 45 C.F.R. 164.514(b)(2) (the HIPAA Privacy Rule) are removed; and the submitting institution has no actual knowledge that the remaining information could be used alone or in combination with other information to identify the subject of the data.

And vet...

This is a requirement, by NIH, that grantee institutions stipulate to a known falsehood in order to receive funding for genomic-scale sequencing.

no actual knowledge that the remaining information could be used alone or in combination with other information to identify the subject of the data.





Re-identification is not the problem.

The Delusion of Deidentification Is

This is the second post in Bill of Health's symposium on the Law, Ethics, and Science of Re-Identification Demonstrations. We'll have more contributions throughout the week, and

This IS a Big Deal

Big Deal

The (almost trivial) possibilities for data reidentification, using freely available Big Data, means that NIH should (must?) entertain a serious re-analysis of the entire issue of data privacy.

Big Deal

The (almost trivial) possibilities for data reidentification, using freely available Big Data, means that NIH should (must?) entertain a serious re-analysis of the entire issue of data privacy.

We must either develop a new approach that recognizes the possibility (likelihood!) of reidentification, or ...

Big Deal

The (almost trivial) possibilities for data reidentification, using freely available Big Data, means that NIH should (must?) entertain a serious re-analysis of the entire issue of data privacy.

We must either develop a new approach that recognizes the possibility (likelihood!) of reidentification, or ...

prepare to abandon data sharing completely.

Possible New Approach

Excerpts:

For decades an essential principle of privacy laws around the world has been to put individuals in control by letting them decide whether, how, and by whom their personal information may be processed. In the Internet age, this laudable ideal has often morphed into a formulaic system of "notice and consent." In the era of big data, however, when much of data's value is in secondary uses that may have been unimagined when the data was collected, such a mechanism to ensure privacy is no longer suitable.


Possible New Approach

Excerpts:

We envision a very different privacy framework for the bigdata age, one focused less on individual consent at the time of collection and more on holding data users accountable for what they do. ... This spurs creative reuses of the data, while at the same time it ensures that sufficient measures are taken to see that individuals are not hurt.



Possible New Approach

Excerpts:

Running a formal big-data use assessment correctly and implementing its findings accurately offers tangible benefits to data users: they will be free to pursue secondary uses of personal data in many instances without having to go back to individuals to get their explicit consent. On the other hand, sloppy assessments or poor implementation of safeguards will expose data users to legal liability, and regulatory actions such as mandates, fines, and perhaps even criminal prosecution. Data-user accountability only works when it has teeth.



Possible New Approach

Excerpts:

Shifting the burden of responsibility from the public to the users of data makes sense for a number of reasons. They know much more than anybody else, and certainly more than consumers or regulators, about how they intend to use the data. By conducting the assessment themselves (or hiring experts to do it) they will avoid the problem of revealing confidential business strategies to outsiders. Perhaps most important, the data users reap most of the benefits of secondary use, so it's only fair to hold them accountable for their actions and place the burden for this review on them.

