CS-5630 / CS-6630 Uisualization for Data Science Sets and Text

Alexander Lex <u>alex@sci.utah.edu</u>





I HAVE A HARD TIME KEEPING TRACK OF WHICH CONTACTS USE WHICH CHAT SYSTEMS.

Design Workshop

item1 : A item2 : A item3 : A, B item4 : A, C item5 : A, B, C item6 : B item7 : B, C item8 : C

• • •



Venn diagram

LETTER

The banana (*Musa acuminata*) genome and the evolution of monocotyledonous plants

Angélique D'Hont¹*, France Denoeud^{2,3,4}*, Jean-Marc Aury², Franc-Christophe Baurens¹, Françoise Carreel^{1,5}, Olivier Garsmeur¹, Benjamin Noel², Stéphanie Bocs¹, Gaëtan Droc¹, Mathieu Rouard⁶, Corinne Da Silva², Kamel Jabbari^{2,3,4}, Céline Cardi¹, Julie Poulain², Marlène Souquet¹, Karine Labadie², Cyril Jourda¹, Juliette Lengellé¹, Marguerite Rodier-Goud¹, Adriana Alberti², Maria Bernard², Margot Correa², Saravanaraj Ayyampalayam⁷, Michael R. Mckain⁷, Jim Leebens-Mack⁷, Diane Burgess⁸, Mike Freeling⁸, Didier Mbéguié-A-Mbéguié⁹, Matthieu Chabannes⁵, Thomas Wicker¹⁰, Olivier Panaud¹¹, Jose Barbosa¹¹, Eva Hribova¹², Pat Heslop-Harrison¹³, Rémy Habas⁵, Ronan Rivallan¹, Philippe Francois¹, Claire Poiron¹, Andrzej Kilian¹⁴, Dheema Burthia¹, Christophe Jenny¹, Frédéric Bakry¹, Spencer Brown¹⁵, Valentin Guignon^{1,6}, Gert Kema¹⁶, Miguel Dita¹⁹, Cees Waalwijk¹⁶, Steeve Joseph¹, Anne Dievart¹, Olivier Jaillon^{2,3,4}, Julie Leclercq¹, Xavier Argout¹, Eric Lyons¹⁷, Ana Almeida⁸, Mouna Jeridi¹, Jaroslav Dolezel¹², Nicolas Roux⁶, Ange-Marie Risterucci¹, Jean Weissenbach^{2,3,4}

Bananas (*Musa* spp.), including dessert and cooking types, are giant sequence errors. The assembly consisted of 24,425 contigs and 7,513 perennial monocotyledonous herbs of the order Zingiberales, a scaffolds with a total length of 472.2 Mb, which represented 90% of sister group to the well-studied Poales, which include cereals. the estimated DH-Pahang genome size. Ninety per cent of the Bananas are vital for food security in many tropical and subtropical assembly was in 647 scaffolds, and the N50 (the scaffold size above countries and the most popular fruit in industrialized countries¹. which 50% of the total length of the sequence assembly can be found) was 1.3 Mb (Supplementary Text and Supplementary Tables 1–3). We The Musa domestication process started some 7,000 years ago in anchored 70% of the assembly (332 Mb) along the 11 Musa linkage Southeast Asia. It involved hybridizations between diverse species groups of the Pahang genetic map. This corresponded to 258 scaffolds and subspecies, fostered by human migrations², and selection of and included 98.0% of the scaffolds larger than 1 Mb and 92% of the diploid and triploid seedless, parthenocarpic hybrids thereafter widely dispersed by vegetative propagation. Half of the current annotated genes (Supplementary Text, Supplementary Table 4 and Supplementary Fig. 1). production relies on somaclones derived from a single triploid genotype (Cavendish)¹. Pests and diseases have gradually become We identified 36,542 protein-coding gene models in the Musa adapted, representing an imminent danger for global banana progenome (Supplementary Tables 1 and 5). A total of 235 microRNAs duction^{3,4}. Here we describe the draft sequence of the 523-megabase from 37 families were identified, including only one of the eight microRNA gene (MIR) families found so far solely in Poaceae⁸ genome of a Musa acuminata doubled-haploid genotype, providing a crucial stepping-stone for genetic improvement of banana. We (Supplementary Tables 6 and 7). Viral sequences related to the banana streak virus (BSV) dsDNA detected three rounds of whole-genome duplications in the Musa lineage, independently of those previously described in the Poales plant pararetrovirus were found to be integrated in the Pahang lineage and the one we detected in the Arecales lineage. This first genome, with 24 loci spanning 10 chromosomes (Supplementary monocotyledon high-continuity whole-genome sequence reported Text and Supplementary Fig. 2). They belonged to a badnavirus phylogenetic group that differed from the endogenous BSV species outside Poales represents an essential bridge for comparative genome analysis in plants. As such, it clarifies commelinid-(eBSV) found in *M. balbisiana*⁹ and most of them formed a new

Nature 2012



Phoenix dactylifera 28,889 / 19,027

Figure 4 | Six-way Venn diagram showing the distribution of shared gene families (sequence clusters) among *M. acuminata, P. dactylifera, Arabidopsis thaliana, Oryza sativa, Sorghum bicolor* and *Brachypodium distachyon* genomes. Numbers of clusters are provided in the intersections. The total number of sequences for each species is provided under the species name (total number of sequences/total number of clustered sequences).

(

Sorghum bico 34,496 / 27,39

Phoenix dactylifera 28,889 / 19,027

Dryza sativa 10,612 / 27,049

Arabidopsis thaliana 27,169 / 21,950





[Wiles et al., BMC Systems Biology]

[D'Hont et al., Nature, 2012]







Element ID

Name

Lisa

Bart

Homer

Mr. Burns

Charact School, School, Power P Evil, Pov

What are some questions we'd like to ask?

Sets	Attribute(s)
eristics	Age
Female	8
Male	10
Plant, Male	40
wer Plant, Male	90

Element ID	
Name	Charac
_isa	School
Bart	School
Homer	Power
Mr. Burns	Evil, Po

Don't always try to show all individuals
What is the biggest intersection?
Which sets make up an intersection?
How big is an intersection?
Does it work for more than four sets?

Sets	Attribute(s)
naracteristics	Age
hool, Female	8
hool, Male	10
ower Plant, Male	40
vil, Power Plant, Male	90

Design Workshop

- work in groups
- get to know the data (5 mins)
- create three (rapid!) prototypes (3x10 mins)
- Write up your two favorites (15 mins) in google docs

- We'll show you some of our solutions next time!
- Upload to "Bonus" Canvas Dropbox by 4pm

Uenn and Euler Diagrams

Venn vs Euler Euler Diagram Shows logical relations May omit empty intersections



Venn Diagram Shows all possible logical relations between sets (even if empty)



Venn Diagrams

Venn diagrams for many sets are hard

of intersections is 2ⁿ



https://en.wikipedia.org/wiki/Venn_diagram



Area-Proportional Euler Diagrams

- Problem with Venn: size doesn't correspond to the data.
- Creating area-proportional Euler diagrams is hard.
- Layout criteria:
 - simple curves (circles are best)
 - makes it easy to identify which sets are participating in intersection
 - Gestalt-principle: good continuation
 - area proportional



[Alsallakh 2015]

Compare Simple vs Complex Shape

Complex



Simple





[created with EulerAPE]

>< 19

22



[created with EulerAPE]

Venn-Euler Pros/Cons

Pros Familiar Intuitive Work well for 2-4 sets

Cons Don't work well for more than 4 sets Area proportional hard to do Not well suited to show attributes

Relationships for specific Items



No Duplicate Nodes **Complex Shapes** Notice the Nesting



Duplicate Nodes Simple Shapes

[Riche 2010]









Sets on top of a fixed layout



https://www.youtube.com/watch?v=Ju2hSThmPWA

Sets on top of a fixed layout

LineSets



Kelp Diagrams



Node-Link Techniques

Treat sets as nodes Connect to elements that are in set



Instant Messaging

. A.



.







Showing Pairwise Overlap

- Shows fairways overlap of sets
- Doesn't show higher-order overlaps
- Very scalable
- Can't show attributes

Co-Mutations of genes



Pairwise + Interaction



Set Matrices: OnSet

- Set membership for each item shown in matrix
- Comparisons can be made using AND or OR operations
- Good for many sets and few items







Linear Diagrams







Fig. 2. Visualizing sets: Euler diagrams.





	School	B? Hour	Duff for	Buil	male	Po phont	Ase	
Lisa	1/1/1							
Bart	11111				1nm			
Homer			Inn		1114	Intra		1
marge		11/11			1000	-0000		+
Marge	1							-
maggie								+
Barney			Im		1111			-
Mr Burn				1/10	lin	rim	1	
Mo			Im		Im	1	There	



Radial Sets

- Sets are segments on a "circle"
- Relationships are encoded as ribbons
- Size of segments encodes size of sets
- Histograms in segments show degrees







1	1 - 10
	11 - 2
mine	21 - 3
-	31
3.77, 2.	41 -
[.: : .]	512



Evil







[InfoVis'14]

UpSet Visualizing Intersecting Sets



			Perpagan
Release Date	Average Rating	Times Watchec	₽2,000 -
1,950 2,000	2 4	0 2,000	1,500 -
			1,000 -
			500 -
			1.0 1.5 2.0 2.5
			Average Rating
			+ Scatterplot
			Element Queries
			433 🖸 170 名
	<u> </u>		+
	n		
			Query Filters
			O 🖋 Range Avera
·	· ·		Minimum = 4 Maximum = 5
. I			
⊢I	⊢I	⊢	Name ‡
HD	HDH		Query Results
		Ι	
⊢ □ ਮ			Name
⊢–––-DH		₩ □ I	Toy Story (1995)
⊢DH	$\vdash - \Box - \downarrow$		Sense and Sensibilit
⊢⊡H	⊢[]	ll I	Persuasion (1995)
⊢ []	н		(1995)
			(1990)



age Rating

Name	Release Date	Average Rating	Tim Wat
Toy Story (1995)	1995	4.15	207
Sense and Sensibility (1995)	1995	4.03	835
Persuasion (1995)	1995	4.06	179
City of Lost Children, The (1995)	1995	4.06	403

Contains

÷





Sel Vis Croals

1. Efficient visual encoding



2. Creating complex slices of a dataset

3. Visualize attributes







[Movie Lens Dataset]



Visualizing Intersections

Visualizing Properties

Attribute Details

	Seven (Se7en) (1995)		

Element List & Queries





Visualizing Intersections








Universal Set







PLOEEING ARTIPULES



Massisisistic intersection? attribute in an intersection ots















Which is the biggest intersection? Sort By: Cardinality









Aggregation



Are many items shared between two sets? **Aggregate By: Degree**









Are many items shared between two sets? **Aggregate By: Degree**

Sum of children









How are the elements of 'B' distributed? **Aggregate By: Set**







How are the elements of 'B' distributed? **Aggregate By: Set**







How are the elements of 'B' distributed? **Aggregate By: Set**







Queries







Element Visualizations

No visualizations configured. Click + button to add a new visua

Scatterplot



Element Queries

No queries. Click + button to add a new query.

Query Filters

+

No active query.

Query Results

No active query.





Altribulices



3

How do documentaries compare to adventure movies?





How do documentaries compare to adventure movies?



1996

1996

White Squall (1996)

Muppet Treasure Island (1996)



R-Version: UpSetR

Developed at HMS Some design adaptions







The Banana Chart Redesigned







DEGIGINE



		-
uff Fan	Bluettin Values.	
	DITI3	-
0	0 172	
0	0 []]2	
0	• -	
	0 0	
•	0	
0	• •	
•	0 0	
il	Power Plant	
	0 1114	
0	0 []2	
0	• □1	-
•	• □	-
	0	
•	0 0	H
0	0	-
-		Ħ
here a stall		

Other Options



http://setviz.net

Design Critique



https://goo.gl/IDRXDI

http://mariandoerk.de/edgemaps/demo/



Text and Document Visualization

Slides adapted from Hendrik Strobelt

Text / Language

Features of Text as representation language abstract, general extremely expressive different across population groups (countries, accents, religions,...) linear perception semi-structured (content: grammar, words, sentences, paragraphs,..; appearance: typography, calligraphy,..)

Why Visualize Text?

Worldwide Corporate Data Growth



Design and Text

- Typography:
 - typefaces (serif, sans-serif, **bold**, *italic*)
 - point size (10pt, 12pt, 24pt, 36pt...) line length (alignment: left, right, justified)
 - vertical: line spacing (leading)
 - horizontal: spaces between groups of letters (tracking)
 - space between pairs of letters (kerning) combining letters to a glyph ligatures

Creating a font type is an art that requires profound design knowledge

 $fi \rightarrow fi$ $f1 \rightarrow f1$





Comic Sans and Higgs Boson

We present updated results on SM Higgs searches based on the data recorded in 2011 at \$\sum s=7 TeV (-4.9 fb⁻¹) and 2012 at \$\sum s=8 TeV (-5.9 fb⁻¹)

Results are preliminary:

- 2012 data recorded until 2 weeks ago
- harsher conditions in 2012 due to ~ x2 larger event pile-up
- new, improved analyses deployed for the first time

H → yy and H→ 41: high-sensitivity at low-m_b; high mass-resolution; pile-up robust
 analyses improved to increase sensitivity → new results from 2011 data
 all the data recorded so far in 2012 here analyzed
 → results are presented here for the fit with the

Other low-mass channels: $H \rightarrow WW^{(n)} \rightarrow IvIv$, $H \rightarrow \pi$, $W/ZH \rightarrow W/Z$ bb:

- □ E_T^{miss} in final state → less robust to pile-up
- worse mass resolution, no signal "peak" in some cases.
- complex mixture of backgrounds
- Inderstanding of the detector performance and backgrounds in advanced, but results not yet mature enough to be presented to
- → 2011 results used here for these channels for the overall combined





Visualization for "Raw" Text

in daily use..

enriched text - hypertext linking (graph navigation)



overview & detail



highlighting semantics

```
oid base64_encode(const uint8_t * data, size_t length, char * dst)
size_t src_idx = 0;
size_t dst_idx = 0;
 or (; (src_idx + 2) < length; src_idx += 3, dst_idx += 4)</pre>
    uint8_t s0 = data[src_idx];
    uint8_t s1 = data[src_idx + 1];
    uint8_t s2 = data[src_idx + 2];
    dst[dst_idx + 0] = charset[(s0 & 0xfc) >> 2];
    dst[dst_idx + 1] = charset[((s0 & 0x03) << 4) | ((s1 & 0xf0) >> 4)];
    dst[dst_idx + 2] = charset[((s1 & 0x0f) << 2) | (s2 & 0xc0) >> 6];
    dst[dst_idx + 3] = charset[(s2 & 0x3f)];
 if (src_idx < length)</pre>
    uint8_t s0 = data[src_idx];
    uint8_t s1 = (src_idx + 1 < length) ? data[src_idx + 1] : 0;</pre>
    dst[dst_idx++] = charset[(s0 & 0xfc) >> 2];
    dst[dst_idx++] = charset[((s0 & 0x03) << 4) | ((s1 & 0xf0) >> 4)];
    if (src_idx + 1 < length)</pre>
        dst[dst_idx++] = charset[((s1 & 0x0f) << 2)];</pre>
```



Visualization for "Raw" Text **Document Lens** Visualizing Search Results



Figure 3: Document Lens with lens pulled toward the user. The resulting truncated pyramid makes text near the lens' edges readable.

> Robertson, George G., and Jock D. Mackinlay The document lens Proceedings of the 6th annual ACM symposium on User interface software and technology. ACM, 1993.

A. Stoffel, H. Strobelt, O. Deussen, D. A. Keim Computer Graphics Forum, volume 31 issue 3 pp. Eurographics Conference on Visualization (EuroVis) 2012 S. Bruckner, S. Miksch, and H. Pfister (Guest Editors)

University

Volume 31 (2012), Number 3

Document Thumbnails with Variable Text Scaling

A. Stoffel and H. Strobelt and O. Deussen and D. A. Keim

of Konstanz, Germany

Abstract

Document reader applications usually offer an overview of the layout for each page as thumbnail view. Reading the text in these becomes impossible when the font size becomes very small. We improve the readability of these thumbnails using a distortion method, which retains a readable font size of interesting text while shrinking less interesting text further. In contrast to existing approaches, our method preserves the global layout of a page and is able to show context around important terms. We evaluate our technique and show application examples

1. Motivation



the occurrence of keywords in the documents. So the user

has to step through all occurrences of the keyword within scrolling the pages



pages . In addition , thumbnails can be useful for retrieval if the users are trying

know [CvDRH99 , DC02]. Due to the small size of text in thumbnails , the highlighting should in addition increase the size of the keywords and their context , at first to make the text better readable and second to allow a simple dis ambiguation of keywords by their context . For instance , it

about "USER" or "USER" inter-face"keyword "USER" would

The technique we present to create the thumbnails is a general distortion technique for document content that high - to a user defined interest The global structure of a page, namely the position of im -

is used that highlights the keywords and their context. Other applications might use a different interest function, for instance a sentiment score could be used to create thumbnails for sentiment analysis

2. Related Worl

Three different techniques are currently used for handling document overview and navigation: abstraction from the document with pixel based representations, thumbnails with different highlighting techniques, and semantic zooming

A common pixel based technique is TileBars [Hea95], which visualizes the length of documents and the distribu tion of search terms within these documents with a rectain gular pixel-based visualization. Byrd [Byr99] combines the scrollbar of the document view with a pixel visualization of



order to access the context of the search terms.

Thumbnails, small version of the document or page, are commonly used for overview and navigation. The spacefilling thumbnail approach of Cockburn et al. [CGA06] avoids scrolling in the overview of a document, by positioning the thumbnails of all pages on a grid on the screen and resizing the thumbnails to fit the window size. Suh et al. [SWRG02] combined the thumbnails with popouts, which highlgiht search terms by rendering them in a readable

Document Thumbnails with Variable Text Scaling
Working with Text

unstructured text





structured data

Structured Text Features

simple counts (bag of words) used for similarity measures



dragon	castle
1	1
0	1

Typical Steps of Processing to derive Text Features

Large collections require pre-processing of text to extract information and align text. Typical steps are:

- cleaning (regular expressions)
- sentence splitting
- change to lower case
- stopword removal (most frequent words in a language)
- stemming <u>demo porter stemmer</u>
- POS tagging (part of speech) <u>demo</u>
- noun chunking
- NER (name entity recognition) demo opencalais
- deep parsing try to "understand" text.

Text features are complicated

Toilet out of order. Please use floor below.

One morning I shot an elephant in my pajamas. How he got in my pajamas, I don't know.

Did you ever hear the story about the blind carpenter who picked up his hammer and saw?

http://en.wikipedia.org/wiki/List_of_linguistic_example_sentences

Text Units Hierarchy



linguistic visualization

document collection visualization

single document visualization

Wordle

Frequency-based words that occur often are large Can vary font type, size, color, etc.

http://www.wordle.net







Wordle vs Tag Cloud



Fig 2: Wordle vs. Tag Cloud of Barack Obama's speech at the Democratic Convention in 2008.

addetion admit afford afghanistan ago agree ahead alive america american americans army auto back benefits breaks bush business businesses Care cars century challenges chance change child children clean clean clean college companies country create cut daughters day days debate decades decent democrats deserve dignity datars dreams drea economic economy education election end energy face failure families family finally find finish fix fundamental fundamentals future generation george give giving good GOVERNMENT grateful great hands hard health hear heard higher home hope idea ideas invest iraq job jobs john judgment kennedy lead leave life lives long longer lost love made make makes making man market mccain measure meet men michelle middle-class military million moment moments money moral nation new night nuclear obligation oil part party past pay people percent plan plans plant politics poverty power president programs. progress promise protect proud provide pursue put ready renewable republicans require responsibility restore ward rise safe security Senator sense set sick sights small stand standards start states stood strength student talk talking tax taxes teachers logy ten things thing threats time today tonight tough troops turn understand united veterans walk washington watch watching whiners woman women WORK worked workers working world years young









Word Tree

Text

if love be rough with you , be rough with love . if love be blind , love cannot hit the mark . if love be blind , it best agrees with night .

WordTree



Search for "if" in romeo & Juliet



The word tree, an interactive visual concordance

M Wattenberg, FB Viégas Visualization and Computer Graphics, IEEE Transactions on 14 (6), 1221-1228

PhraseNets





Many Eyes finds this word relationship in Jane Austen's text:

Her manners were pronounced to be very bad indeed, a mixture of pride and impertinence; she had no conversation, no stile, no taste, no beauty.

Many Eyes creates the word graph:

Frank van Ham, Martin Wattenberg, and Fernanda B. Viegas. Mapping Text with Phrase Nets. IEEE Transactions on Visualization and Computer Graphics 15, 6 (November 2009)

pride-impertinence



Corpora: MDS Approaches

use bag-of-word to project documents w.r.t. text similarity into a landscape (only) one example

Fernando V. Paulovich, Franklina M. B. Toledo, Guilherme P. Telles, Rosane Minghim, and Luis Gustavo Nonato. Semantic Wordification of Document Collections. *Comp. Graph. Forum* 31, 3pt3 (June 2012)



gion (bottom figure).

JigSaw

🛸 Document Cluster View		
Edit View Bookmarks Export Options	;	
Highlight Viewed Documents Filters Image: All Filters	center, faa, air staff, p muslims, saudi, governments	directo professional,member terrorist,unit,state
Group by Filters		
Undo Filters	binalshibh, atta, flight	interview concrel a
Hide Unfiltered		
Clusters Text (Default) Freq Words Unique Words	deputy,staff	
All Documents 	9/11,saudi,deceased	04,fdny agency,office



DocumentCards

DC - pipeline

Interaction:

- caption tooltip
- abstract tooltip
- move to orig. Pos.
- page switch
- term highlighting

Compare Corpora

Compare topics between text collections InfoVis SciVis

Figure 1: Comparison of 495 papers of InfoVis, SciVis, and Siggraph (discrimination threshold = 6, number of topics = 30)

Comparative Exploration of Document Collections: a Visual Analytics Approach (<u>http://ditop.hs8.de</u>) D. Oelke, H. Strobelt, C. Rohrdantz, I. Gurevych, and O. Deussen

Vis for Time-Evolving Document Collections

Marian Dörk, Daniel Gruen, Carey Williamson, and Sheelagh Carpendale. A Visual Backchannel for Large-Scale Events. TVCG: Transactions on Visualization and Computer Graphics (Proceedings Information Visualization 2010

THESE CHARTS SHOW MOVIE CHARACTER INTERACTIONS. THE HORIZONTAL AXIS IS TIME. THE VERTICAL GROUPING OF THE LINES INDICATES WHICH CHARACTERS ARE TOGETHER AT A GNEN TIME.

https://xkcd.com/657/

StoryFlow: Tracking the Evolution of Stories

[Liu 2013]

http://textvis.lnu.se/

Text Visualization Browser

A Visual Survey of Text Visualization Techniques Provided by ISOVIS group

Techniques displayed: 141				
Search:				
Analytic Tasks	All and			1 1
Visualization Tasks			in the second	Ar Anges
Data			have a lead.	
Scurce			Maria and and find the second and th	Non-part Non-part
 		 CBC MO SERVICIAN ∞ → → ∞ → ↓ ∞ → ↓ 	Class BBC Annual Class BBC	

