

INCT-CID



DEXL LAB
EXTREME DATA LAB

Fabio Porto
LNCC – MCTIC – DEXL LAB

Semestre 1 – Setembro 2016

Programa de Pós-Graduação em Ciência da Computação

Escola de Informática & Computação
CEFET/RJ

Agenda

- Introdução
- Os Primeiros Passos
- O Turco
- Aplicações



Laboratório Nacional de Computação Científica (LNCC)



DEX LAB
EXTREME DATA LAB



Petropolis, Rio de Janeiro

LNCC - MCTI

- Graduate Course on Computational Modelling
 - CAPES 6
- LABINFO
 - Genomics, proteomics bacteria and micro-organisms
- INCT –MACC
 - Medicine Supported by Scientific Computing
- INCT- CID
- SINAPAD
 - National System of High Performance Performance
- Thematic Laboratory
 - ACIMA (acima.Incc.br)
 - MARTIN (martin.Incc.br)
 - DEXL (dexl.Incc.br)
 - COMCIDIS (comcidis.Incc.br)
 - HEMOLAB (hemolab.Incc.br)
 - LABINFO (labinfo.Incc.br)

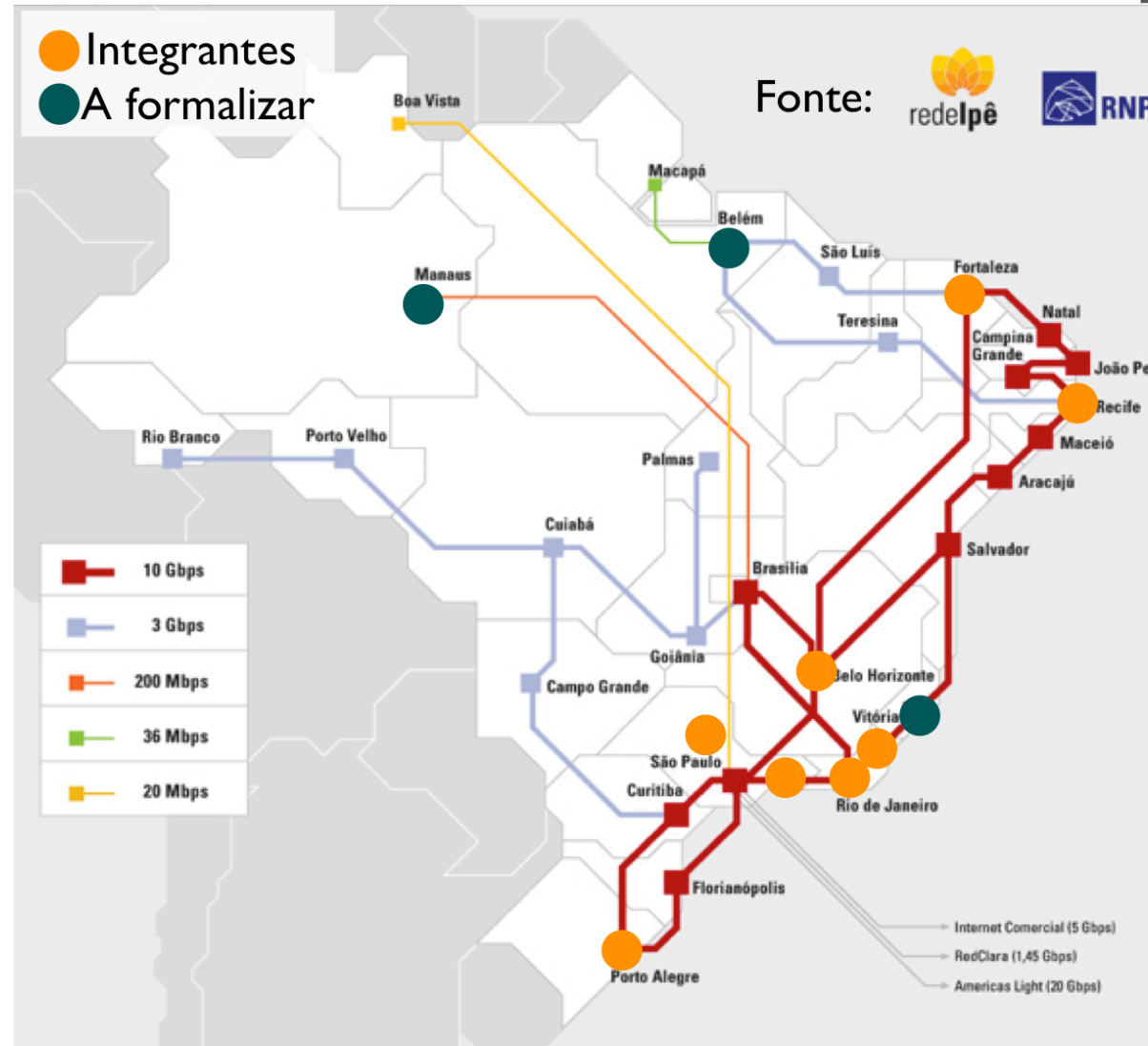


SINAPAD – National Center for HPC



DEXL LAB
XTREME DATA LAB

- Organized in CENAPADS:
 - Universities
 - Research centers
 - Different architectures:
 - Shared Disks
 - Shared Memory
 - GPUs



Atos-BULL Supercomputer



SANTOS DUMONT CONFIGURATION



Machine Partition	Blades	CPU		Total Cores CPU
		Processors/Blade	Cores/processor	
B710 CPU's	504	2	12	12.096
B715+GPU	396	1	12	4.752
B715+PHI	108	1	12	1.296
B710+B715 CPU's	1008		12	18.144
B710cp+B715cpuPhi	612		12	13.392
MESCLA node			16	250
GPU/PHI	3,7			
	Blades	GPU		Total Cores GPU or PHI
		Processors/Blade	Cores/processor	
NVIDIA K40	396	1	2880	1.140.480
INTEL 7120X	108	1	61	6.588

Grupo de Pesquisa em Ciência de Dados



- Artur Ziviani
- Eduardo Ogasawara (CEFET-RJ)
- Fabio Porto
- Kary Ocana
- Luiz Gadelha
- Andre Salles

Ciência de Dados – Uma nova Disciplina



- Mudança de Visão Histórica para Predição do Futuro
 - Predições rasas e profundas
- Inferências baseadas em dados
- Modelos Baseadas em dados
- Foco
 - Grande volume
 - Treinamento
 - Tratamento
 - Interpretação Humana

INCT - Ciência de Dados (CID)



- Objetivo
 - Estruturar a nova área
 - Formação de Recursos Humanos
 - Pesquisa e Desenvolvimento
 - Transferência de Tecnologia para Indústria e Sociedade
 - Como?
 - Através da Integração de disciplinas com foco em dados
 - Bancos de Dados
 - Análise de Dados
 - Modelagem de Dados
 - Metodologias
 - Dirigida à hipóteses
 - Análise estatística
 - Tecnologias
 - Frameworks Big Data: Spark, Giraph, Tensor Flow,
 - Linguagens de programação: R, python, scala
 - Sistemas: NoSQL, Greenplum, ...
 - Machine Learning: scikit-learn, Apache Mahout,..

INCT - CID



- Parcerias
 - FIOCRUZ
 - Observatório Nacional
 - Comitê Olímpico Brasileiro
 - Prefeituras: Curitiba, Rio de Janeiro
- Empresas
 - EMC - Centro de Pesquisa
 - IBM Research
 - Philips Research
- Conselho Internacional
 - NYU
 - INRIA
 - UCSB
 - Boston University – Data Science Institute

Uma breve história da Ciência de Dados



- *R.A. Fisher, 1935, The design of Experiments*
Correlation does not imply causation
- *Hans P. Luhn, 1958, Business Intelligent Systems, IBM*
 - *Automatic method to provide current awareness services to scientists and engineers*
- *J.W. Tukey, 1977, Exploratory Data Analysis*
 - *.. For seeing what the data can tell us beyond the formal modeling and hypothesis testing .. – inspirou o desenvolvimento do pacote estatístico S*
- *H. Dresner, 1989, Business Intelligence*
 - *The database view on data analytics*

Uma breve história da Ciência de Dados



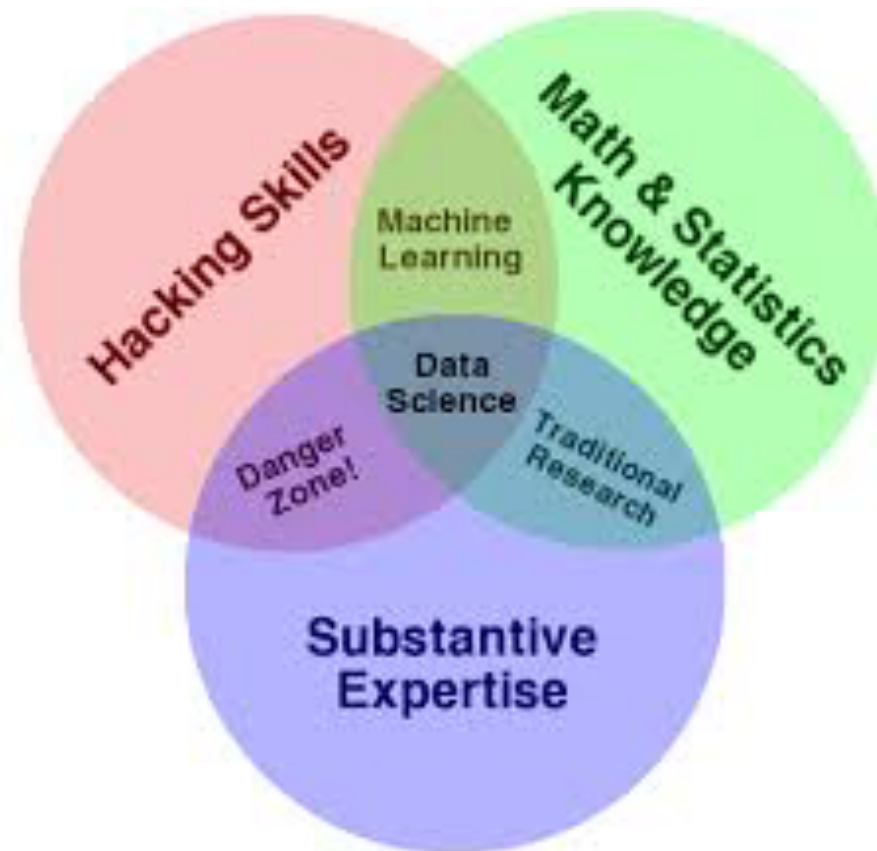
DEXL LAB
EXTREME DATA LAB

- *T. Michel, 1977, The Machine Learning book*
- *Google, 1996, Prototype Search Engine*
- *Jim Gray et. al, 2007, The Fourth Paradigm*
 - *Experimental data science*
- *A. Halevy, P. Norvig, 2009, The Unreasonable Effectiveness of Data*
 - *From deep models to harnessing of data volumes*
- *Exponential growth in data volume, 2010, The data deluge*

Data Science - Habilidades



DEXL LAB
EXTREME DATA LAB

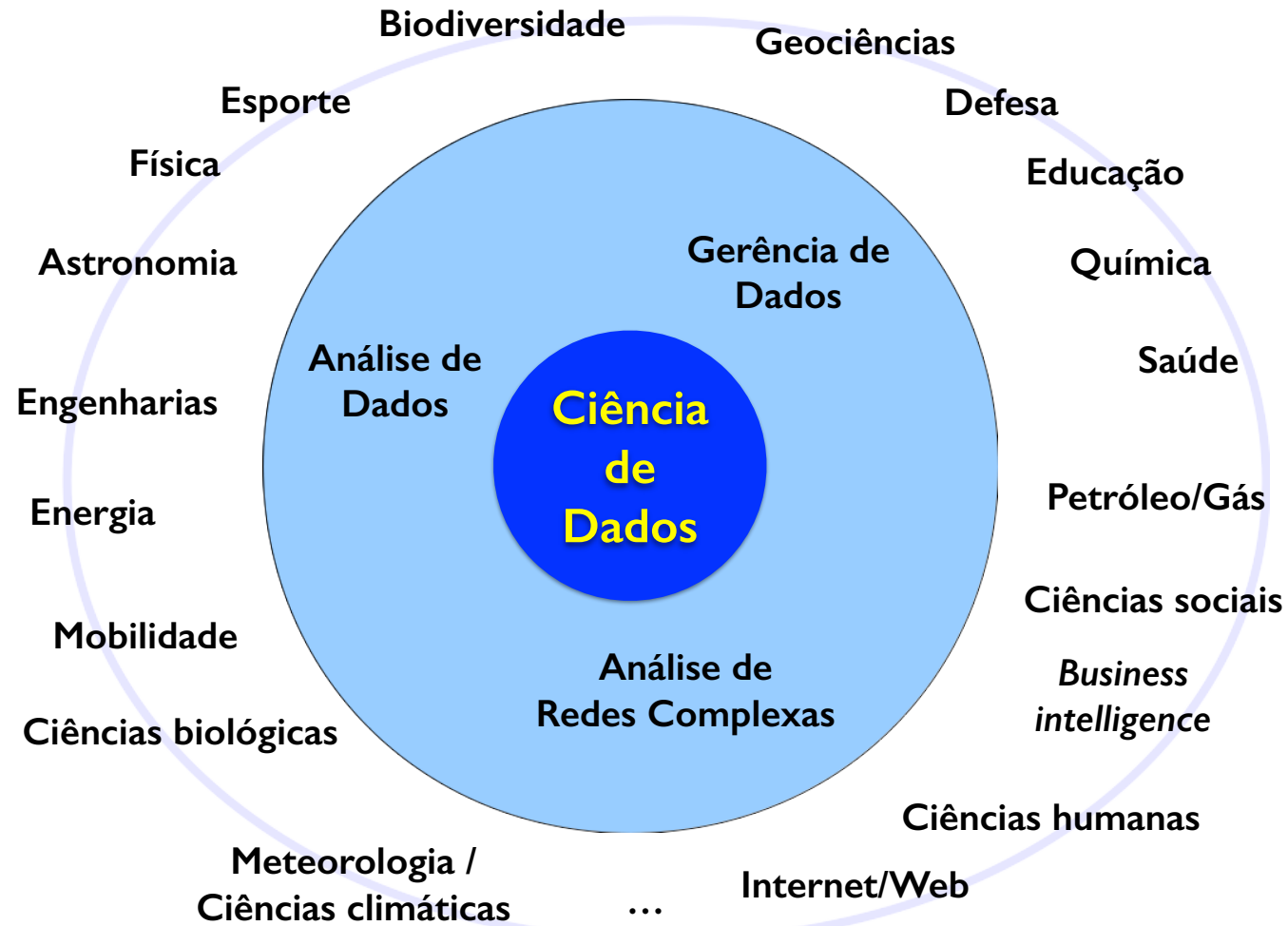


Fonte: <http://drewconway.com/>

Organização do INCT-CID



DEXLAB
EXTREME DATA LAB



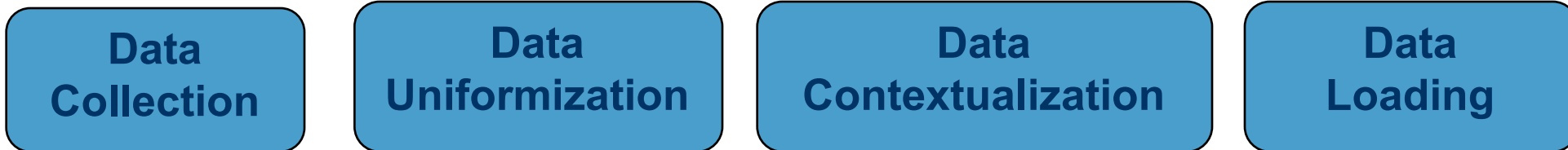
CS + X



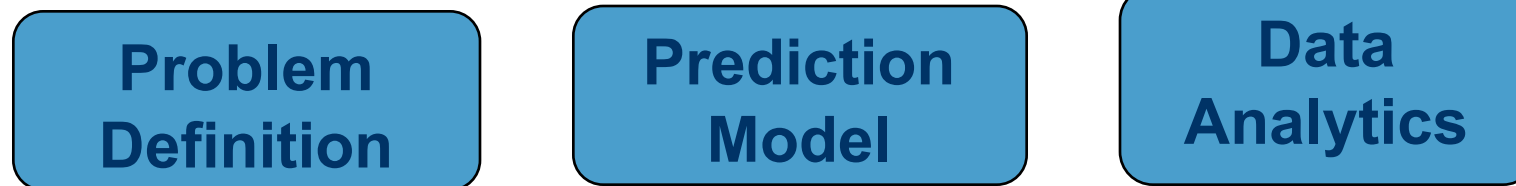
DEXL LAB
EXTREME DATA LAB

- Ciência da Computação está na base do processo científico em todas disciplinas
 - CS + astronomia
 - CS + biologia
 - CS + meteorologia
 - CS + esporte
 - CS + sociologia
- you name it !!
- Fundamental na ciência de dados
- *The Fourth Paradigma: Data Intensive Scientific Discovery, Jim Gray et al.*

Processo de CID



Data Management

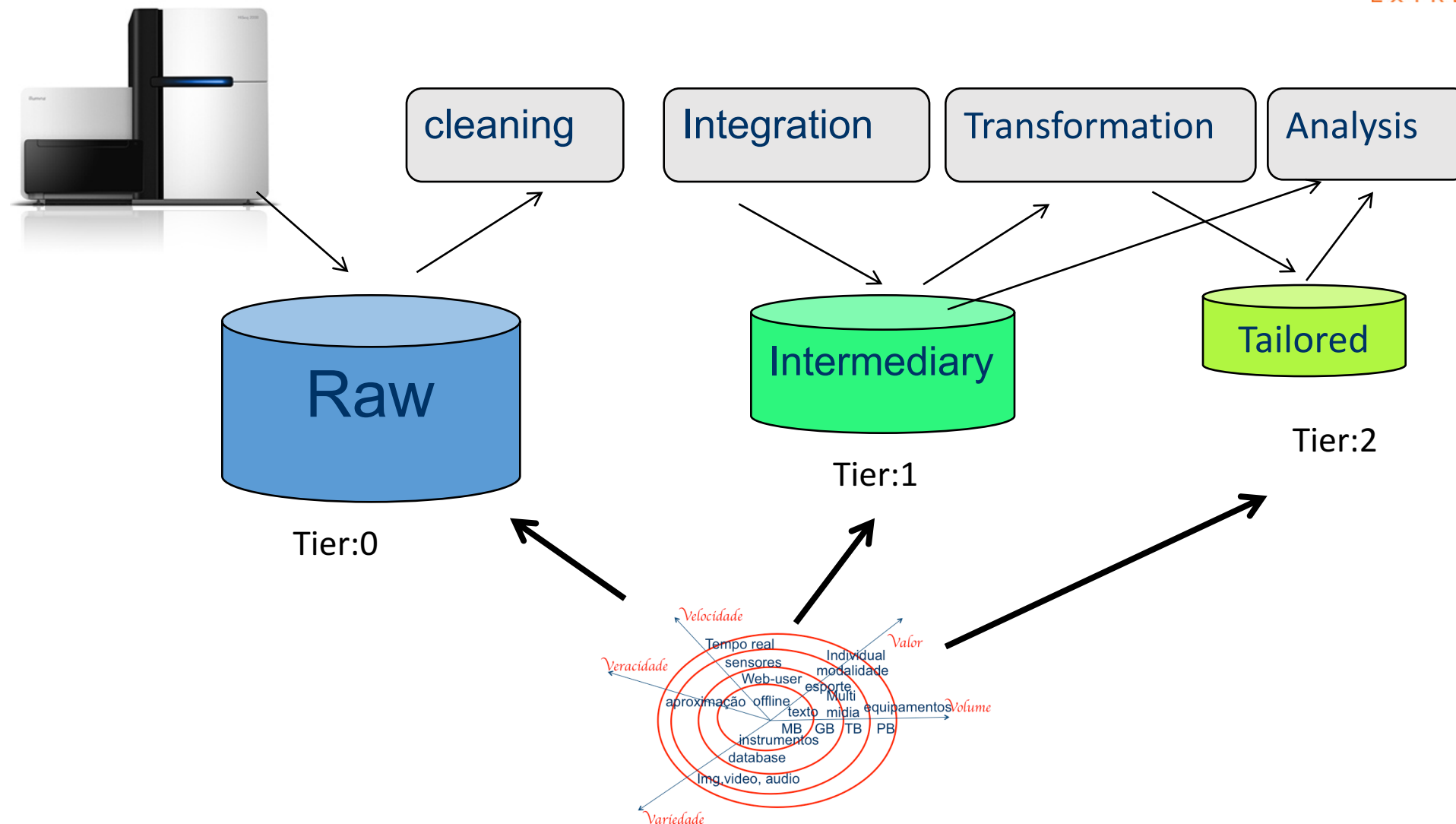


Data Analysis

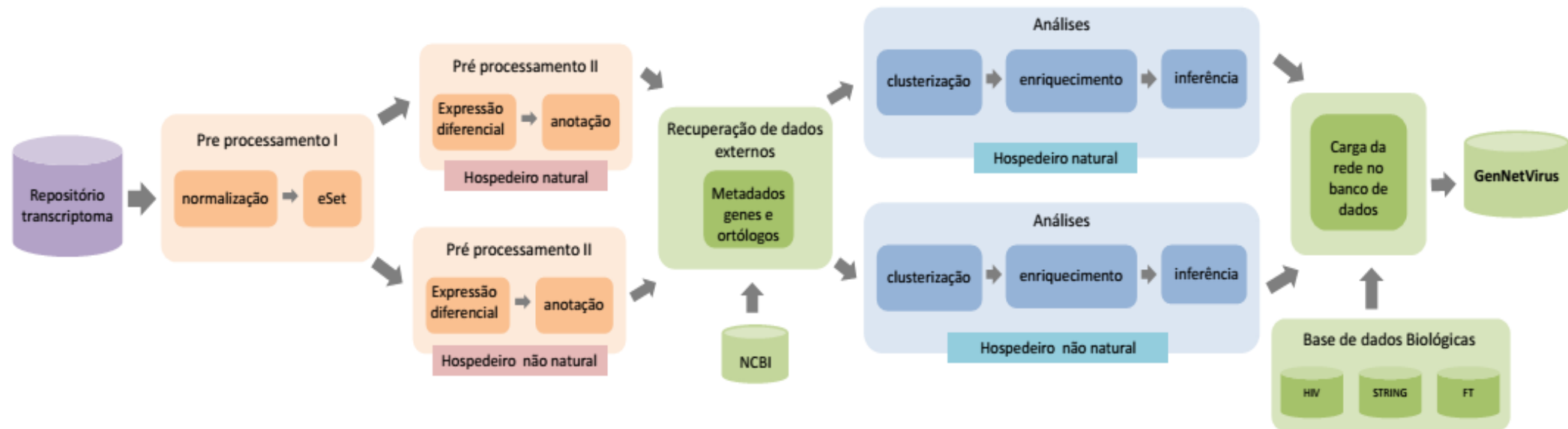
De Dados Crus à Interpretação



DEXL LAB
EXTREME DATA LAB



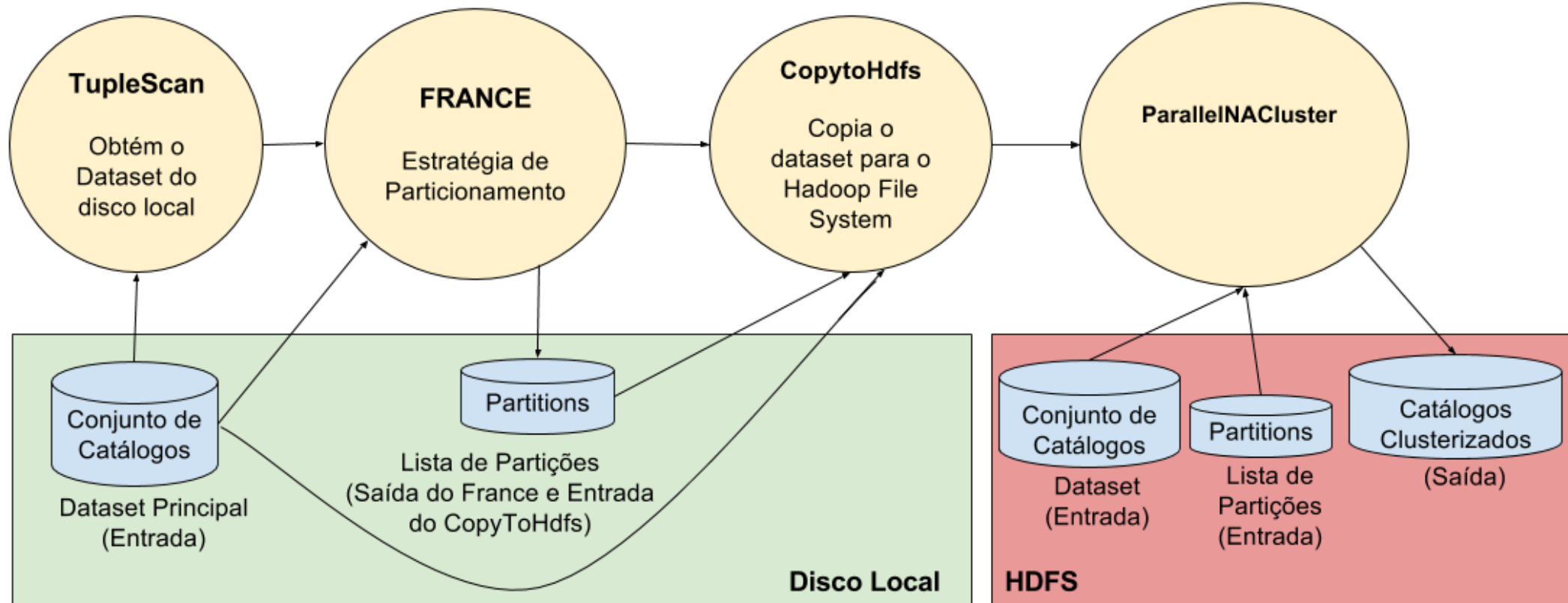
Um dataflow para carga de dados consolidado – GenNetVirus (R. L. Costa)



NaCluster – V. Pires



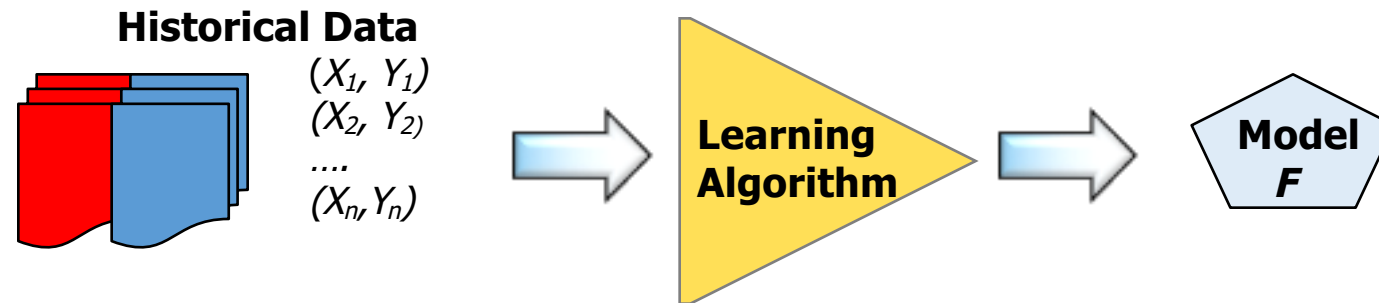
DEXL LAB
EXTREME DATA LAB



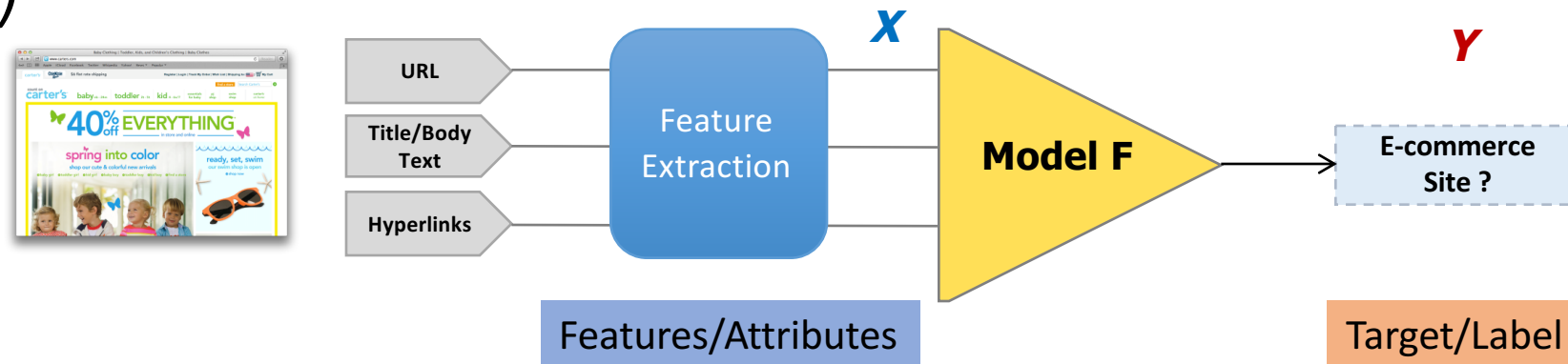


Supervised Learning

- **Training:** Given training examples $\{(X_i, Y_i)\}$ where X_i is the feature vector and Y_i the target variable, learn a function F to best fit the training data (i.e., $Y_i \approx F(X_i)$ for all i)



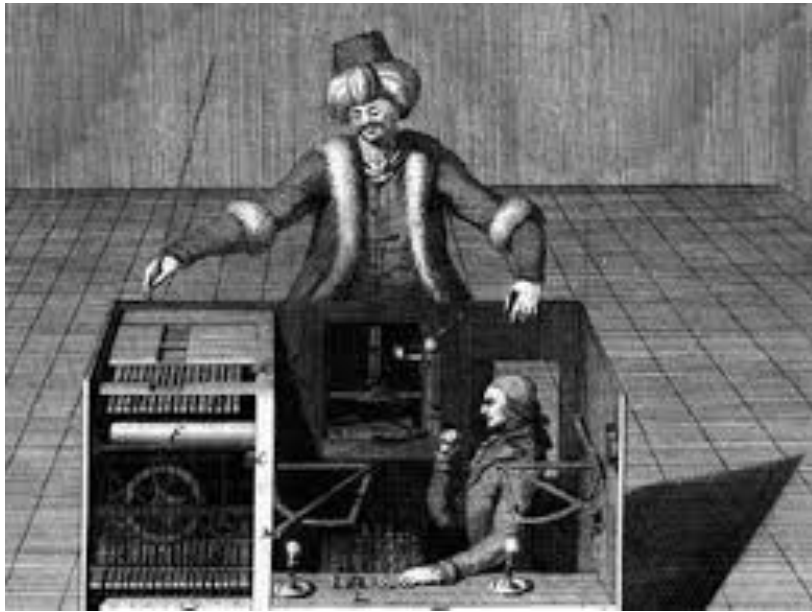
- **Prediction:** Given a new sample X with unknown Y , predict Y using $F(X)$



Human Data Analytics



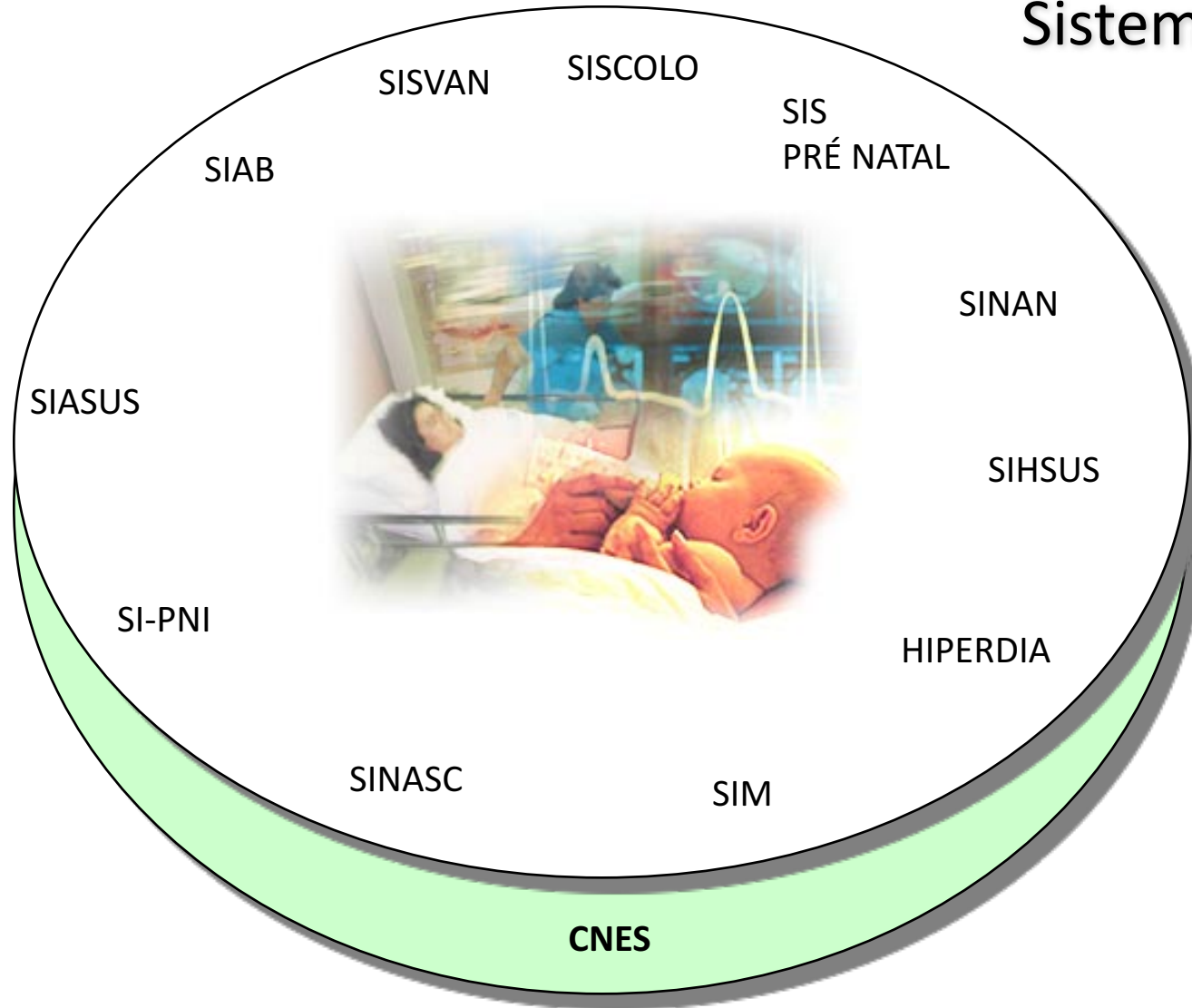
DEX LAB
EXTREME DATA LAB



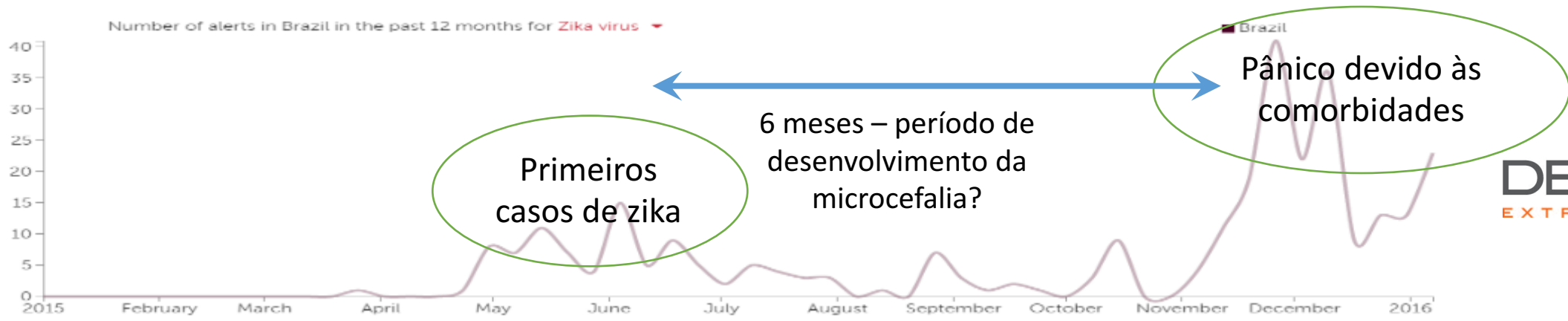
Sistemas de Informação de Saúde



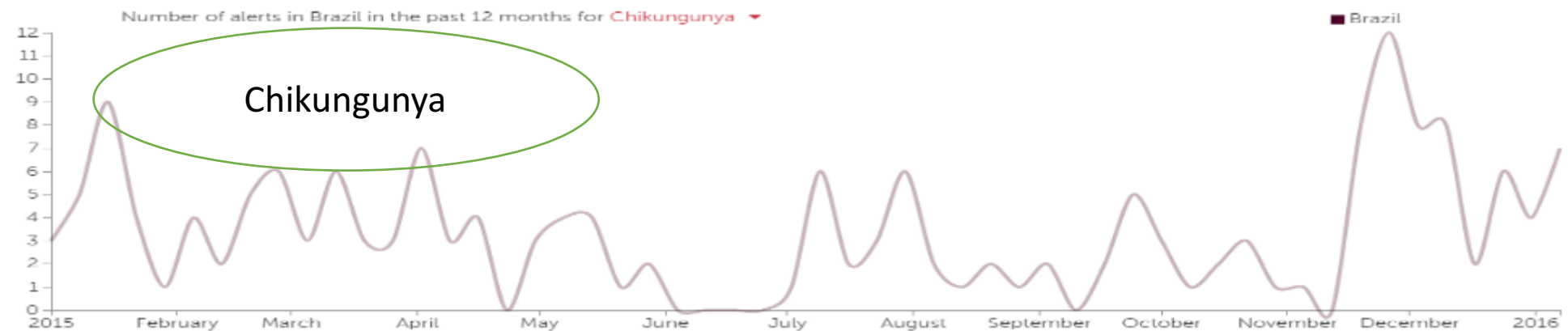
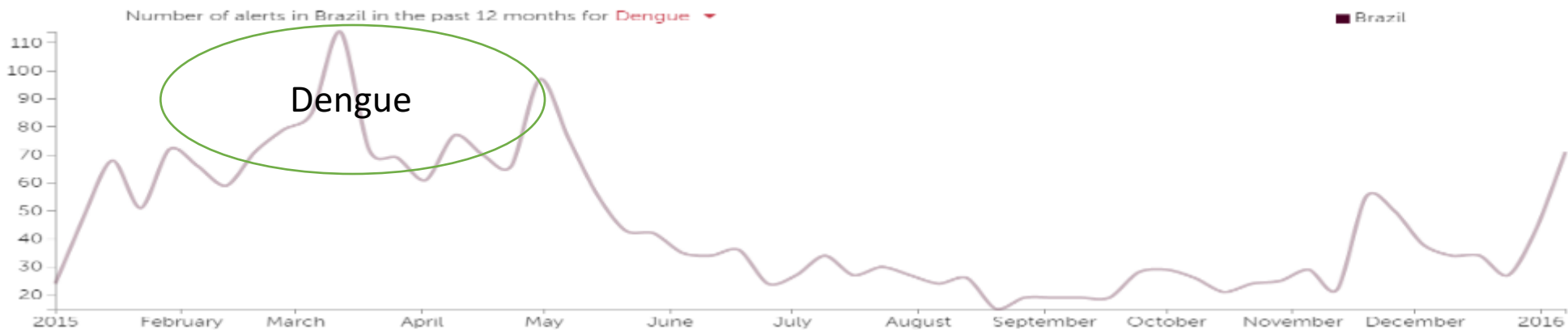
DEXLAB
EXTREME DATA LAB



Registram 'eventos', que ocorrem em pessoas, não os processos de adoecimento e morte
Linkagem de dados, possível mas difícil 24



DEXL LAB
EXTREME DATA LAB



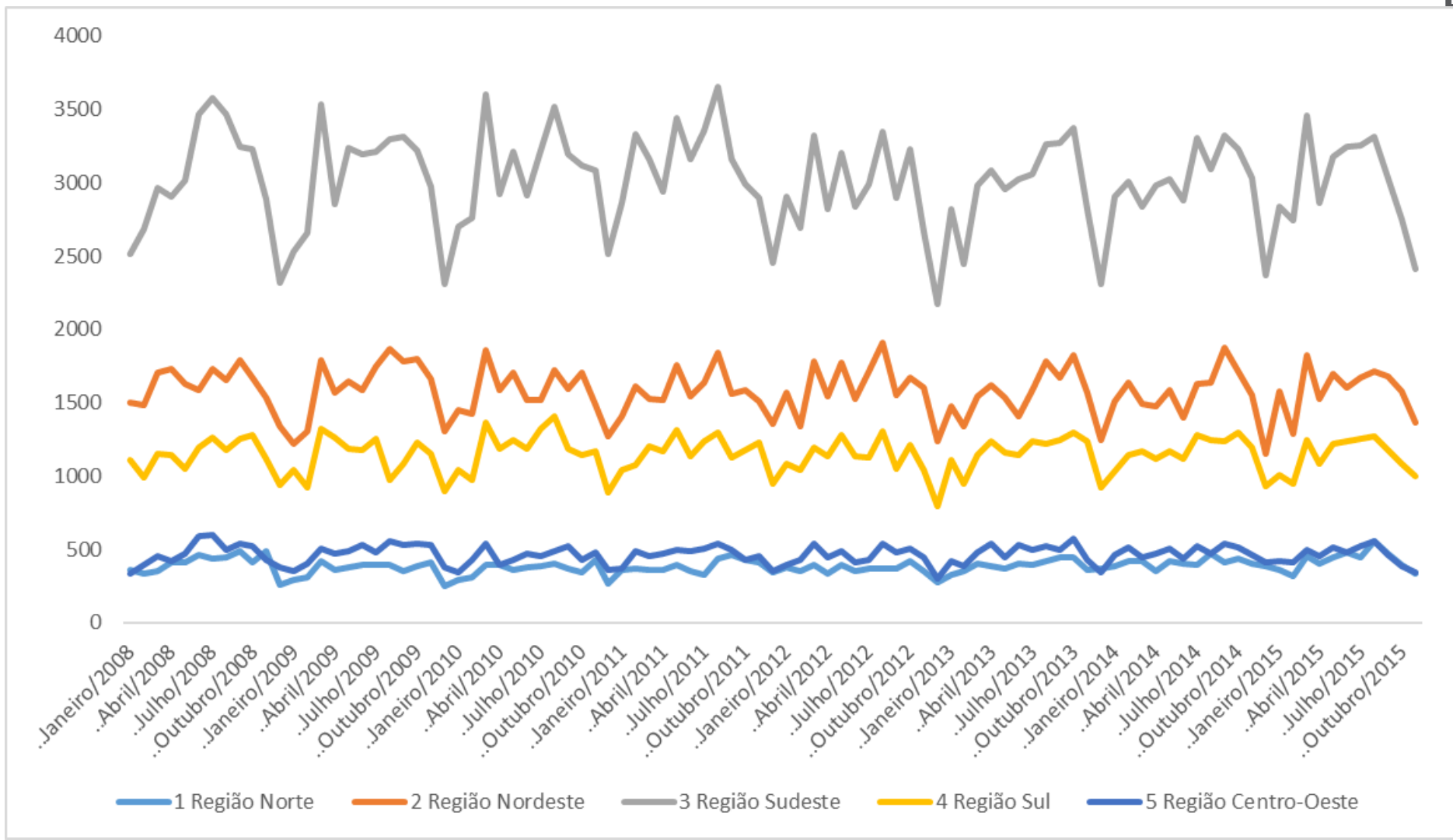
Lista de doenças possivelmente associadas à infecção pelo vírus zika

Doença ou condição	Evento	Sistema	CID-10
INFECÇÕES VIRAIS			
Dengue (dengue clássico)	Notificação	SINAN	A90
Febre hemorrágica devida ao vírus da dengue	Notificação, óbito	SINAN, SIM	A91
Febre de Chikungunya	Notificação	SINAN	A92.0
Doença pelo Zika vírus	Notificação	SINAN	U06*
Doença pelo Zika vírus, não especificada		SINAN	U06.9*
MANIFESTAÇÕES NEUROLÓGICAS ASSOCIADAS A INFECÇÕES VIRAIS (DENGUE, CHIKUNGUNYA E ZIKA)			
Encefalite por vírus transmitidos por mosquitos	Internação hospitalar, óbito	SIH	A83
Outras encefalites por vírus transmitidas por mosquitos		SIH	A83.8
Encefalite não especificada por vírus transmitida por mosquitos		SIH	A83.9
Outras encefalites virais, não classificadas em outra parte	Internação hospitalar	SIH	A85
Encefalite viral não especificada	Internação hospitalar	SIH	A86
Meningite viral	Internação hospitalar	SINAN	A87
Outras meningites virais			A87.8
Meningite viral não especificada			A87.9
Meningite em outras doenças infecciosas e parasitárias classificadas em outra parte	Internação hospitalar, óbito	SIH, SIM	G02
Meningite em doenças virais classificadas em outra parte		SIH, SIM	G02.0
Meningite devida a outras causas e a causas não especificadas	Internação hospitalar, óbito	SIH, SIM	G03
Meningite devida a outras causas especificadas		SIH, SIM	G03.8
Meningite não especificada		SIH, SIM	G03.9
Encefalite, mielite e encefalomielite	Internação hospitalar, óbito	SIH, SIM	G04
Encefalite aguda disseminada (ADEM)		SIH, SIM	G04.0
Outras encefalites, mielites e encefalomielites		SIH, SIM	G04.8
Encefalite, mielite e encefalomielite não especificada		SIH, SIM	G04.9
Encefalite, mielite e encefalomielite em doenças classificadas em outra parte	Internação hospitalar, óbito	SIH, SIM	G05
Encefalite, mielite e encefalomielite em doenças virais classificadas em outra parte		SIH, SIM	G05.1
Síndrome de Guillain-Barré	Internação hospitalar, óbito	SIH, SIM	G61.0
Hemiplegia	Internação hospitalar	SIH	G81
Hemiplegia flácida		SIH	G81.0
Paraplegia e tetraplegia	Internação hospitalar	SIH	G82
Paraplegia flácida		SIH	G82.0
Paraplegia não especificada		SIH	G82.2
Tetraplegia flácida		SIH	G82.3
Síndrome parálitica não especificada	Internação hospitalar	SIH	G83.9
TIPOS E CAUSAS DE ABORTO ASSOCIADOS À INFECÇÃO PELO VÍRUS ZIKA			
Gravidez ectópica	Internação hospitalar, óbito fetal ou da gestante	SIH, SIM-feto	O00
Mola hidatiforme	Internação hospitalar, óbito fetal ou da gestante	SIH, SIM-feto	O01
Outros produtos anormais da concepção	Internação hospitalar, óbito fetal ou da gestante	SIH, SIM-feto	O02
Aborto espontâneo	Internação hospitalar, óbito fetal ou da gestante	SIH, SIM-feto	O03
Aborto por razões médicas e legais	Internação hospitalar óbito fetal ou da gestante	SIH, SIM-feto	O04
Outros tipos de aborto	Internação hospitalar óbito fetal ou da gestante	SIH, SIM-feto	O05
Aborto não especificado	Internação hospitalar óbito fetal ou da gestante	SIH, SIM-feto	O06
Falha de tentativa de aborto	Internação hospitalar óbito fetal ou da gestante	SIH, SIM-feto	O07
Complicações consequentes a aborto e gravidez ectópica ou molar	Internação hospitalar óbito fetal ou da gestante	SIH, SIM-feto	O08
TIPOS DE MALFORMAÇÕES CONGÊNITAS ASSOCIADOS À INFECÇÃO PELO VÍRUS ZIKA			
Anencefalia e malformações similares	Nascimento, óbito fetal, internação do neonato	SINASC, SIM-neonato, SIH	Q00
Microcefalia	Nascimento, óbito fetal, internação do neonato	SINASC, SIM-neonato, SIH	Q02
Hidrocefalia congênita	Nascimento, óbito fetal, internação do neonato	SINASC, SIM-neonato, SIH	Q03
Outra hidrocefalia congênita		SINASC, SIM-neonato, SIH	Q03.8
Hidrocefalia congênita não especificada		SINASC, SIM-neonato, SIH	Q03.9
Outras malformações congênitas do cérebro	Nascimento, óbito fetal, internação do neonato	SINASC, SIM-neonato, SIH	Q04
Outras deformidades por redução do encefalo (incluindo hidranencefalia)		SINASC, SIM-neonato, SIH	Q04.3
Malformação congênita não especificada do encefalo		SINASC, SIM-neonato, SIH	Q04.9
Malformações congênitas das pálpebras, do aparelho lacrimal e da órbita	Nascimento	SINASC	Q10
Anoftalmia, microftalmia e macroftalmia	Nascimento	SINASC	Q11
Malformações congênitas do cristalino	Nascimento	SINASC	Q12
Malformações congênitas da câmara anterior do olho	Nascimento	SINASC	Q13
Malformações congênitas da câmara posterior do olho	Nascimento	SINASC	Q14
Outras malformações congênitas do olho	Nascimento	SINASC	Q15
Malformações congênitas do ouvido causando comprometimento da audição	Nascimento	SINASC	Q16
Outras malformações congênitas da orelha	Nascimento	SINASC	Q17

Evolução das internações por total de malformações em menores de um ano



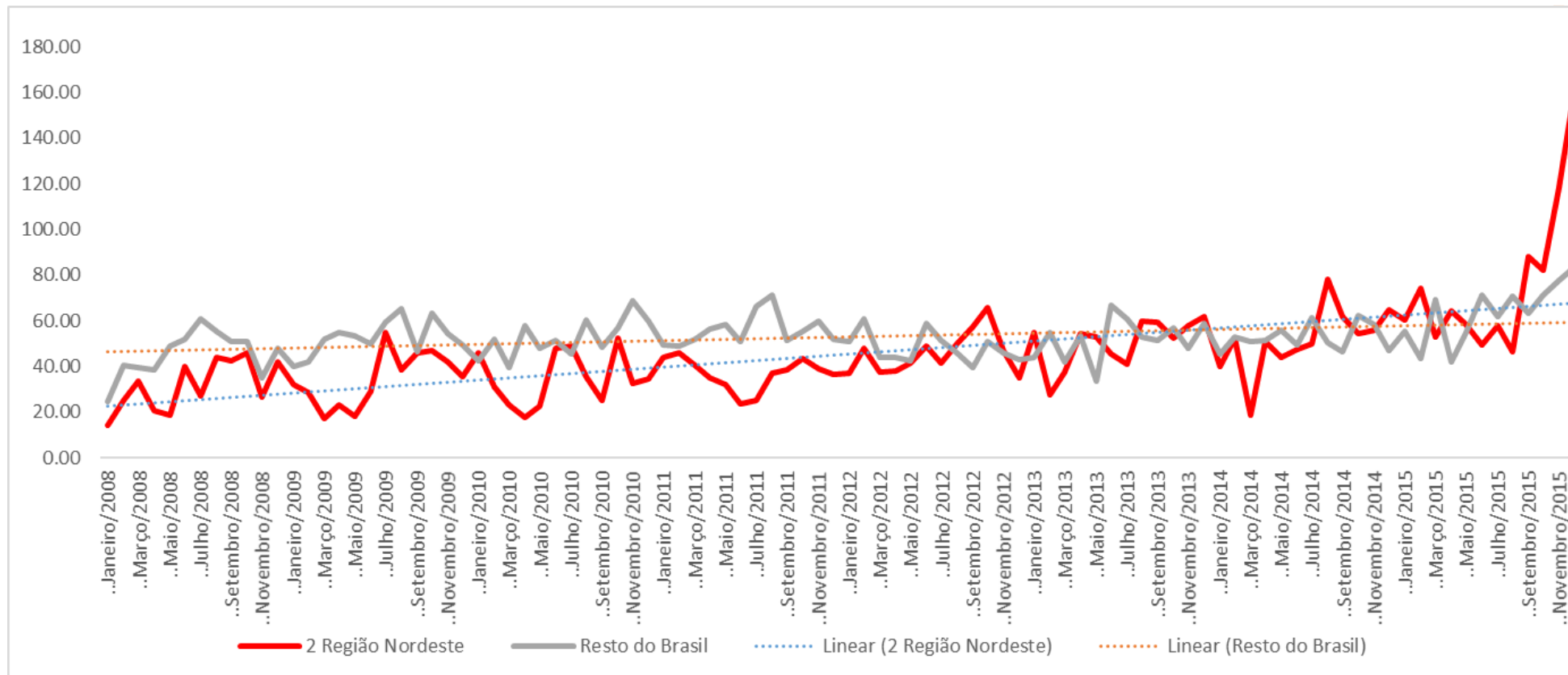
DEXLAB
XTREME DATA LAB



Evolução das interações por malformações do SNC em menores de um ano



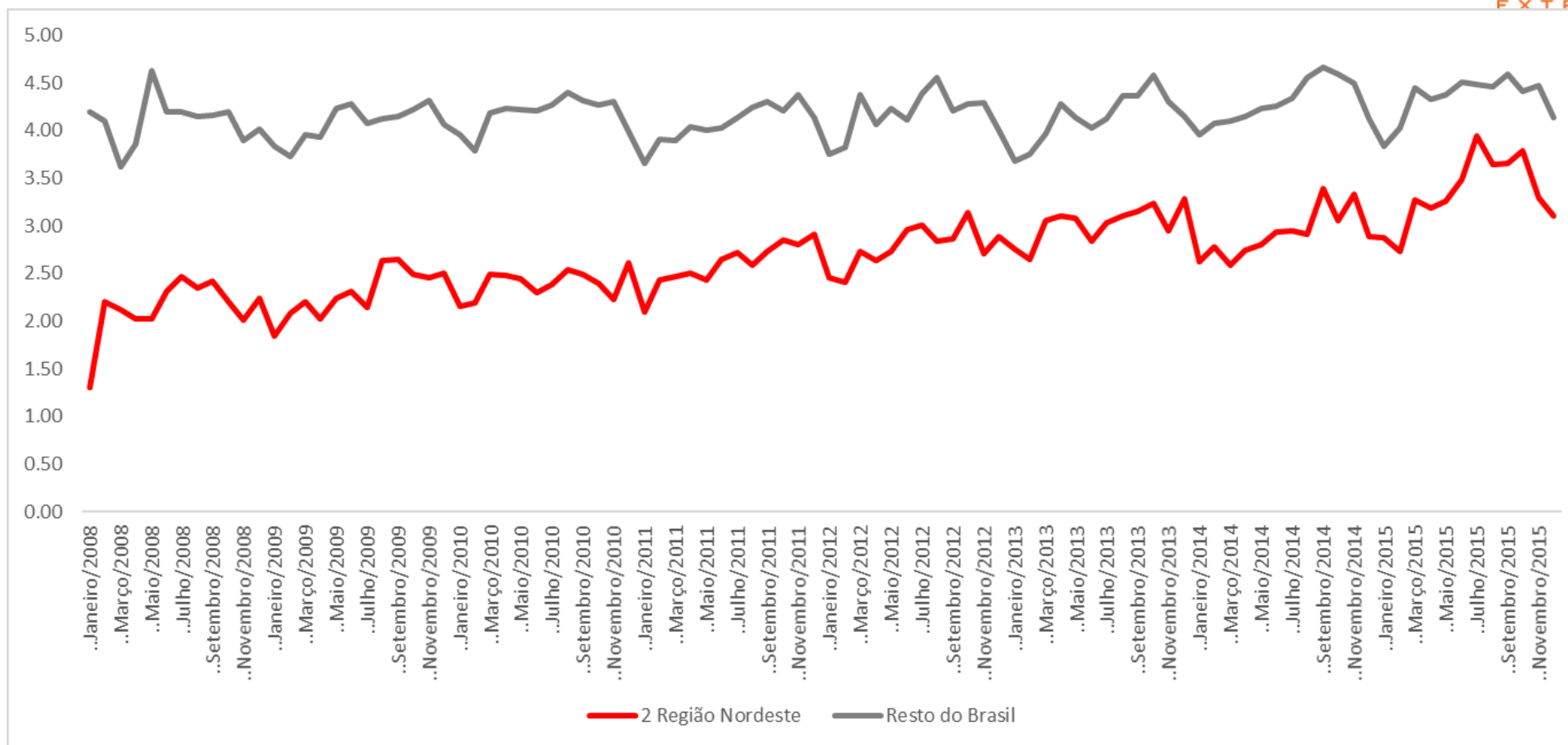
DEXLAB
TREME DATA LAB



Evolução das internações por “outras doenças do SNC”



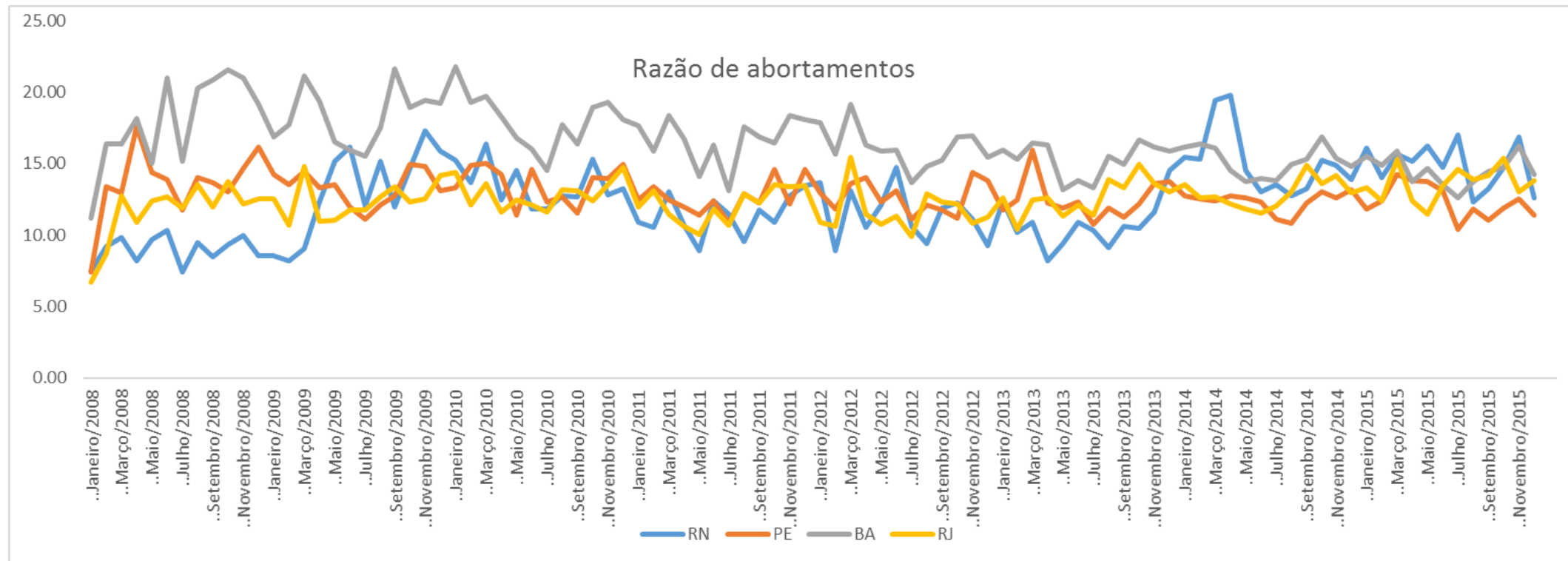
DEXLAB
EXTREME DATA LAB



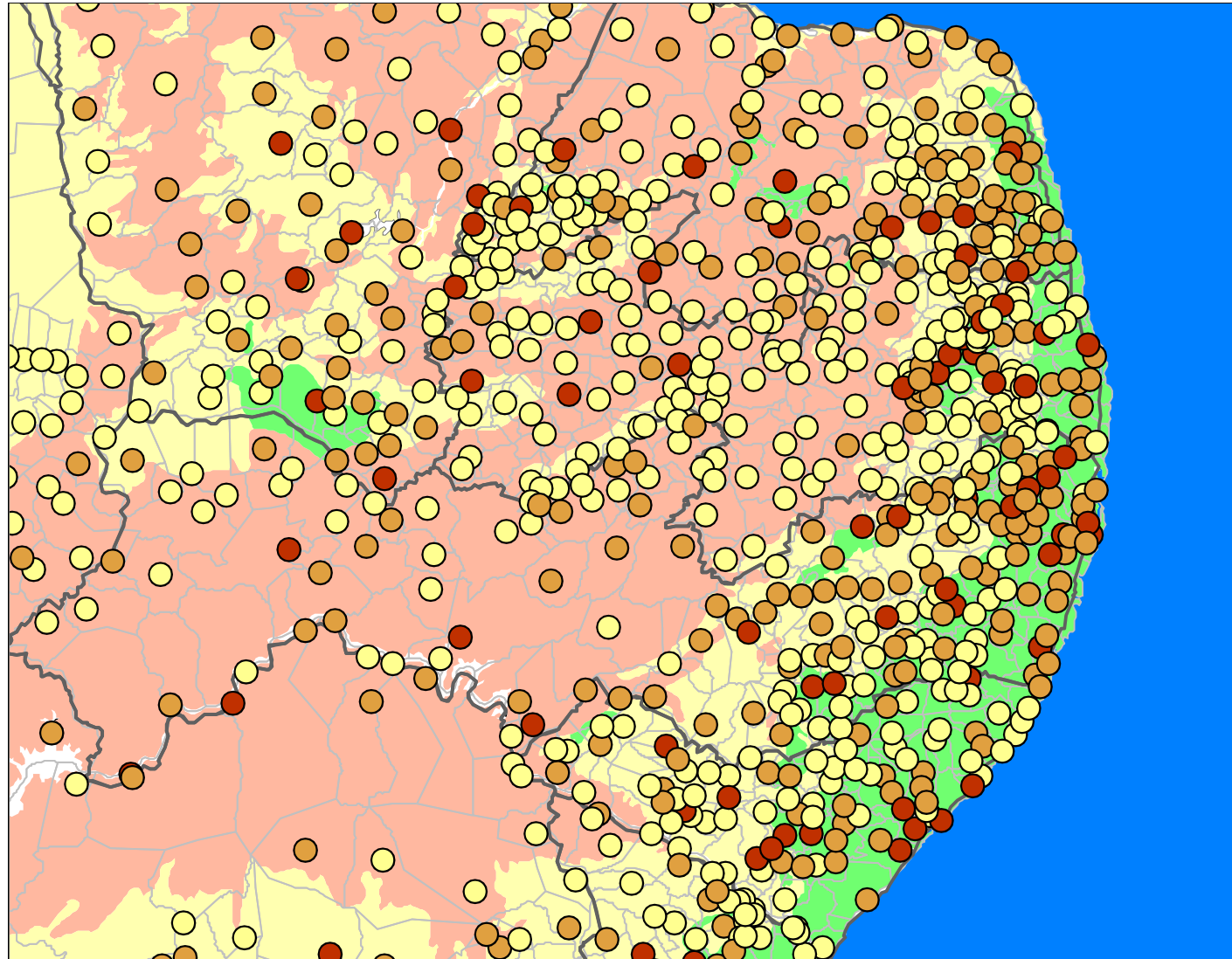
Evolução das internações por complicações relacionadas ao aborto



DEXLAB
EXTREME DATA LAB



Aumento das internações por neuropatias



DEX LAB
EXTREME DATA LAB

Melhorar o indicador

Usar técnicas de detecção de cluster espacial



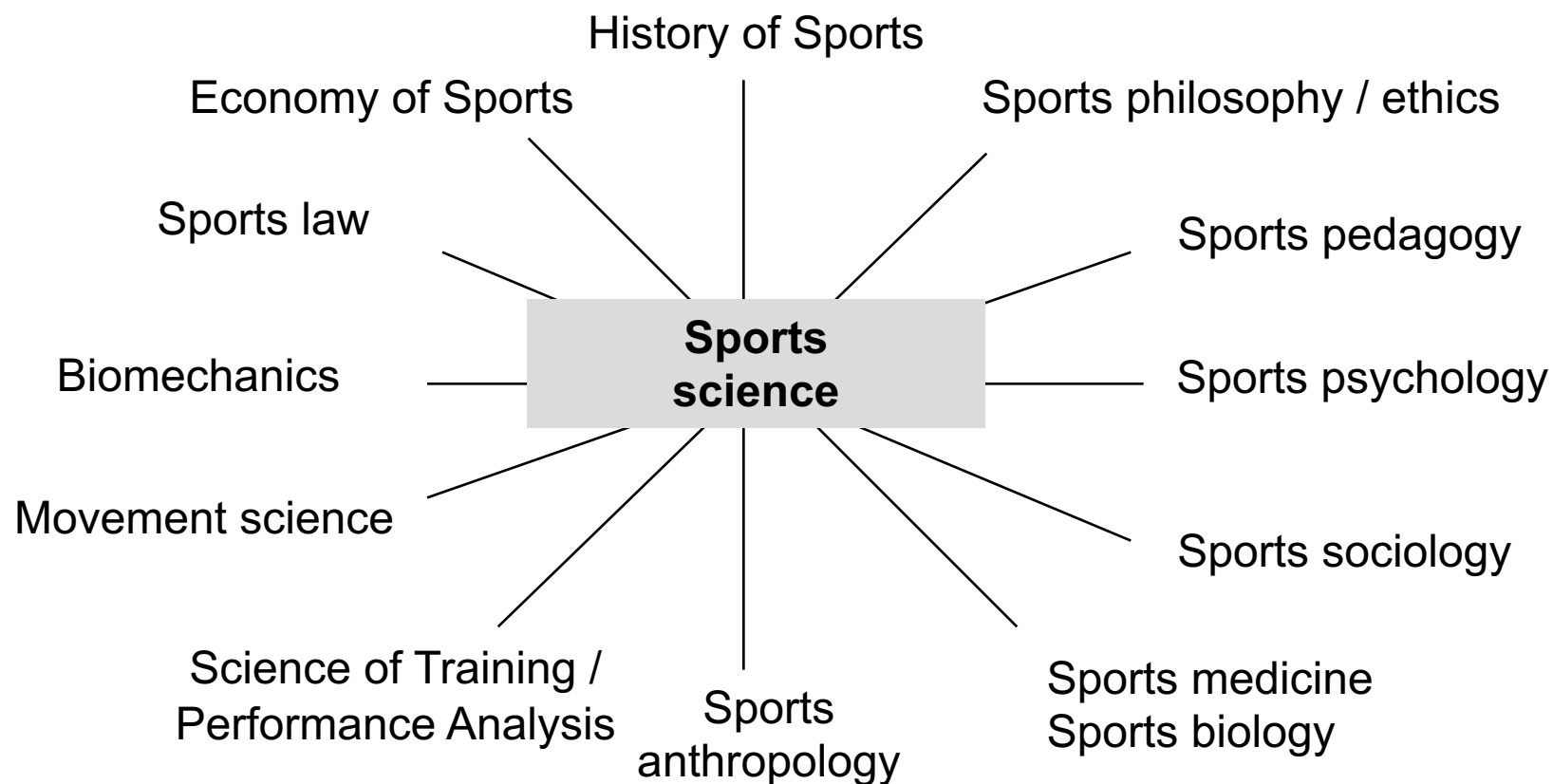
DEXL LAB
EXTREME DATA LAB

Área Exemplo: Esporte de Alto Rendimento



DEXL LAB
EXTREME DATA LAB

Áreas da Ciência do Esporte



Sportomics [Cameron, Bassini 2015]

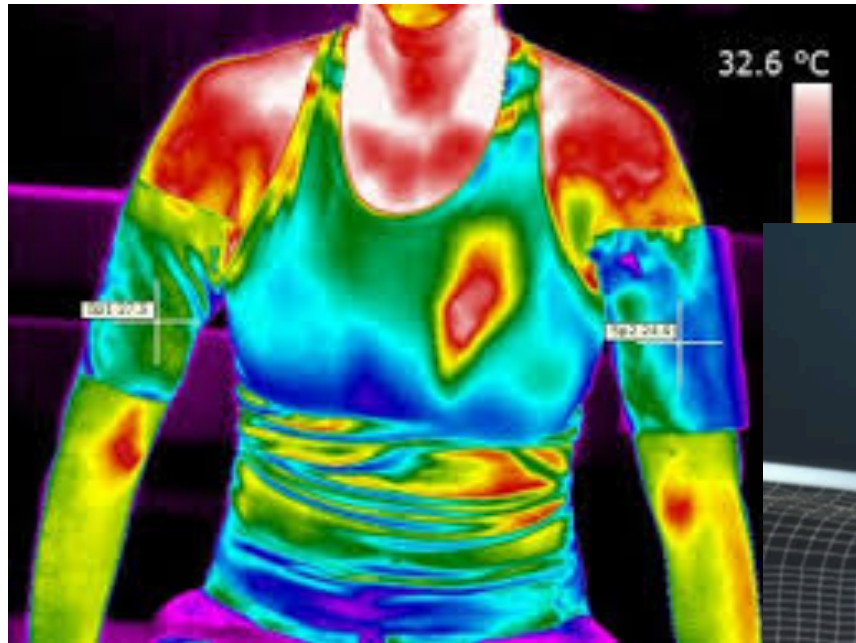


- Adoção da Avaliação integrada de Atletas a partir de diferentes disciplinas
- Características individuais - > fazem a diferença
 - características raras são mais valiosas !!!
- Processo científico baseado em dados

Computação: atleta e sua modalidade



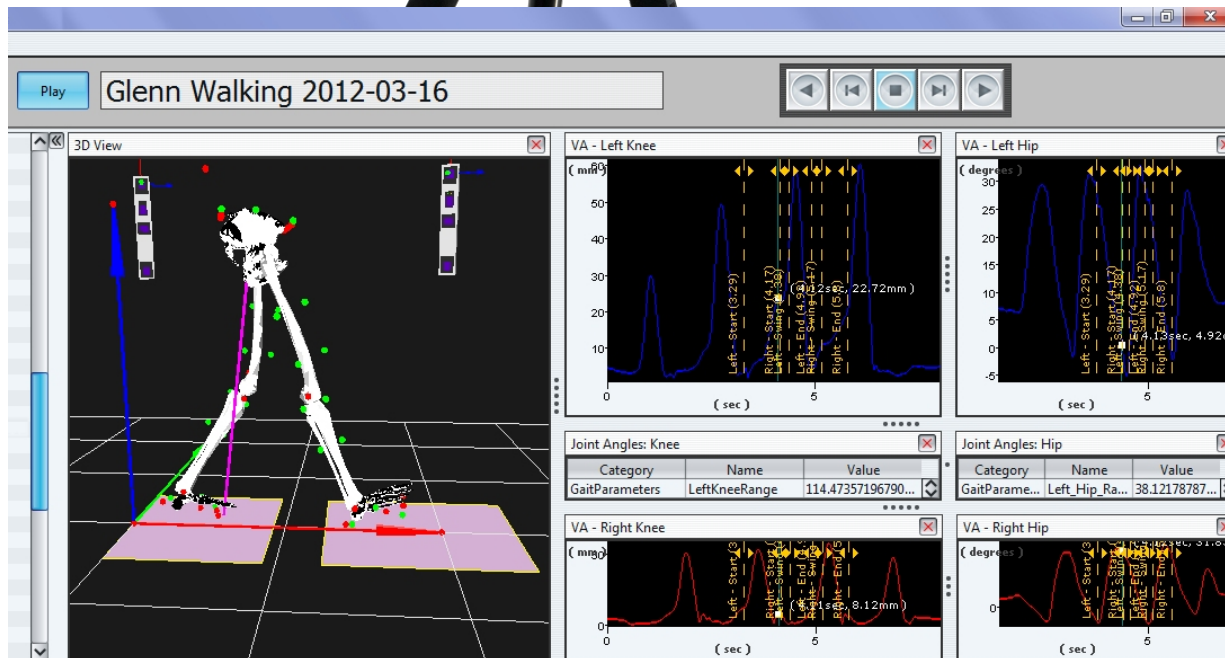
DEXL LAB
EXTREME DATA LAB





DEXL LAB
EXTREME DATA LAB

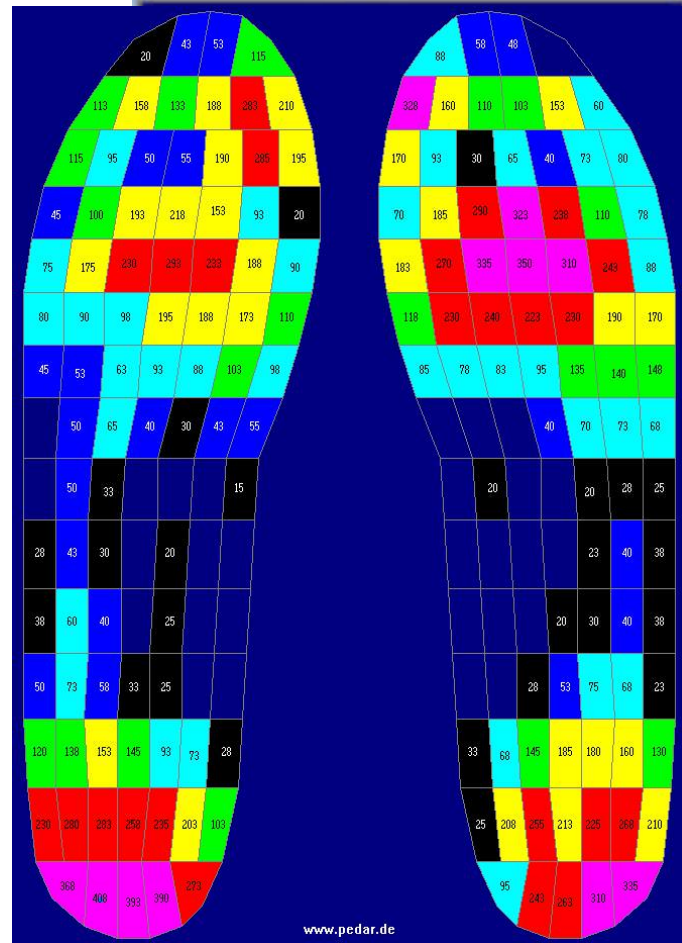
Captura de Dados



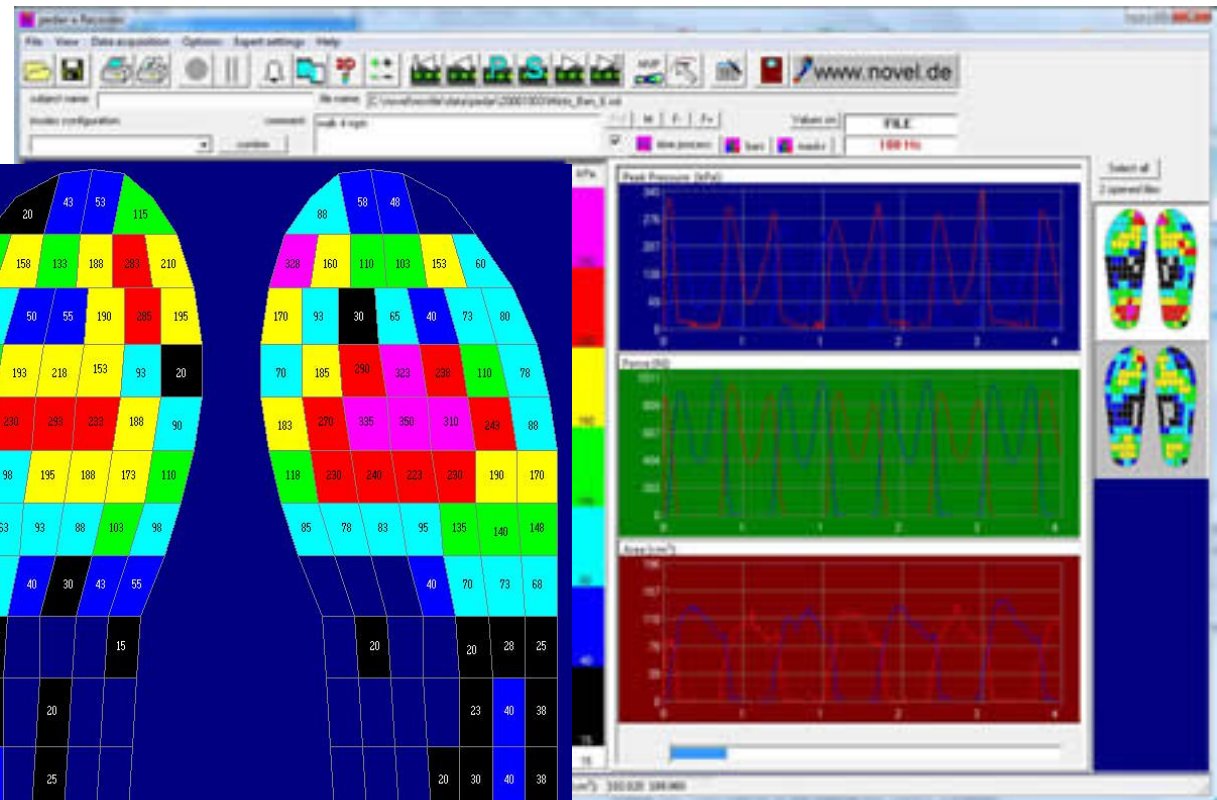
 **DELSYS**[®]
WEARABLE SENSORS
FOR MOVEMENT SCIENCES



novel.de



www.pedar.de



XL LAB
ME DATA LAB



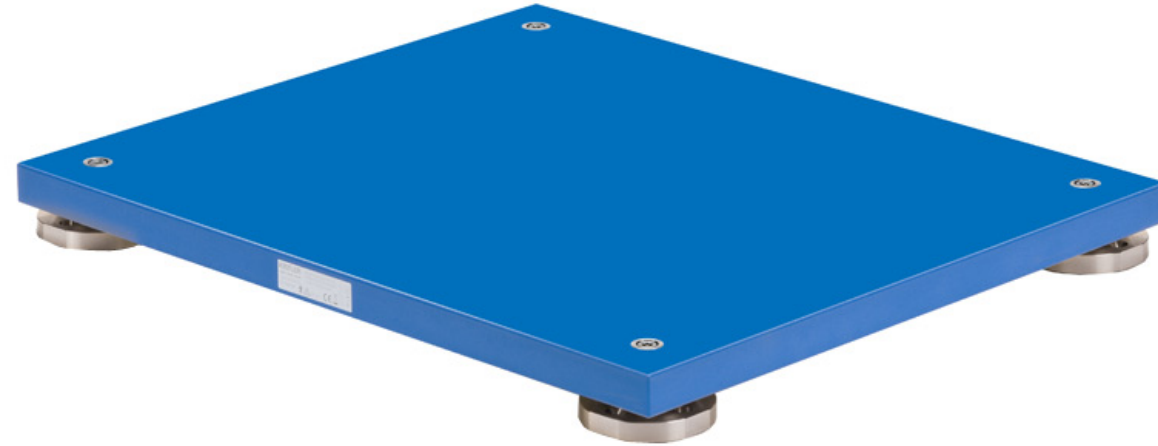
www.novel.de

KISTLER

measure. analyze. innovate.



DFXL LAB
ME DATA LAB



Extração de Sangue



DEX LAB
EXTREME DATA LAB

Espectrometria de Massa



DEXL LAB
EXTREME DATA LAB

[XEVO G2-S TOF]
Identify, quantify, and confirm the broadest range of compounds in the most complex and challenging samples.

StepWave

- Patented T-Wave ion transfer that maximizes sensitivity while minimizing routine maintenance.
- Revolutionary off-axis design for high method robustness with complex sample matrices.

Universal Ion Source

- The most extensive range of ion sources for the broadest range of analytes.
- Tool free design enables ease of use and simplifies routine maintenance.
- Source options are quickly interchangeable and ready-to-use within minutes.

T-Wave Collision Cell

- Enables efficient transfer of ions to the orbitrap mass analyzer.
- Allows rapid MS/MS data acquisition with maximum sensitivity.

Quantof

- Simultaneously delivers UPLC-compatible mass resolution and matrix-tolerant dynamic range, so that the narrowest chromatographic peaks from the most complex sample matrices can be identified and quantified with no compromise in performance.

Dual Stage Reflectron
Improves focusing of ions for high resolution spectra.

High Field Pusher
Increases duty cycle and instrument sensitivity.

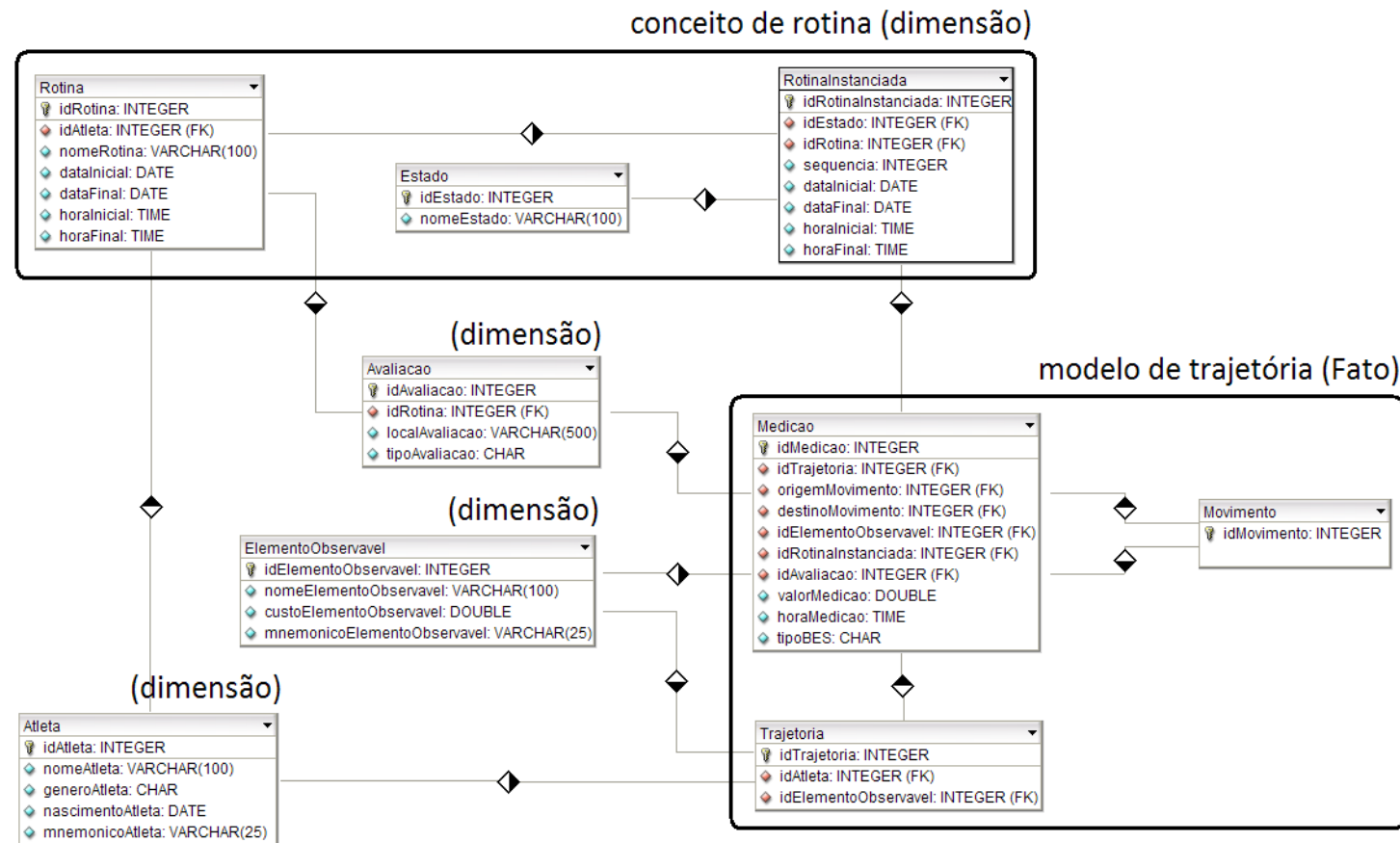
ADC Detector
Provides a linear response over 5 orders of dynamic range.

Waters
THE SCIENCE OF WHAT'S POSSIBLE.™

Modelo Estrela para Análise de Dados



DEXL LAB
EXTREME DATA LAB



Relação: Alvo –Medições



DEXL LAB
EXTREME DATA LAB

Resultado da avaliação do atleta amador: Renato Pinheiro Oliveira

ESPORTE	OBJETIVO	PRAZO	#	DEPTO.	EVIDÊNCIAS	HIPÓTESES	SUGESTÕES
MARATONA	REDUÇÃO DO TEMPO EM 2 MINUTOS	1 ANO	1	BIOQUÍMICA	Análise do nível de CK, GamaGT	Aumentar Treinamento	Implementação de práticas nutricionais suplementares Aumentar treinamento
					Queda da concentração de bicarbonato	Aumentar reserva de bicarbonato	
					Níveis de Creatinina	Melhorar hidratação	
					Níveis de Glicose, Ureia e Creatinina	Rever status nutricional	
					Baixa reserva de glicogênio	Dieta inadequada ou treinamento leve	
			2	NUTRIÇÃO	Baixa ingestão de água	Falta de energia	
					Baixa reposição de energia		
					Baixa ingestão de carboidratos		
					Consumo baixo de vitamina D		
			3	FISIOLOGIA	Análise da composição corporal	Excesso de Gordura	
					Análise do consumo máximo de oxigênio	Baixo consumo de VO2 Baixo Limiar ventilatório 2	
			4	BIOMECÂNICA	Articulação joelho estável	Apto para mais esforço	
						Boa habilidade motora	
						Domínio de força	
			5	COMPORTAMENTO	Índice de auto eficácia	Alta eficácia e Auto Confiança acima da média	
						Índice de Stress/Recuperação	
			6	ANÁLISE DO DESEMP. E SUP. AO TREINO			



DEXL LAB
EXTREME DATA LAB

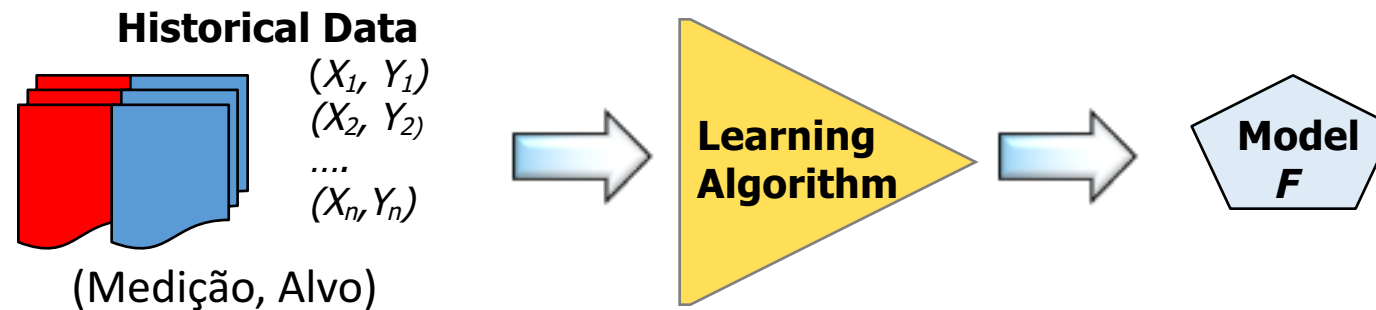
Abordagem

- Identificar atributos relevantes
 - Avaliação estatísticas de correlação entre alvos e medições
- Definir grupos de atributos correlacionados
 - Discretização de atributos (SAX)
 - Determinação de grau de Independência entre Informações Redundante [Wong et al 1976]
 - Uso de algoritmos de clusterização (k-means)
- Estabelecer relações probabilísticas de causalidade
 - Modelo Grafo Probabilístico -> Causal Bayesian Networks
- Identificar elementos a sofrerem interferência

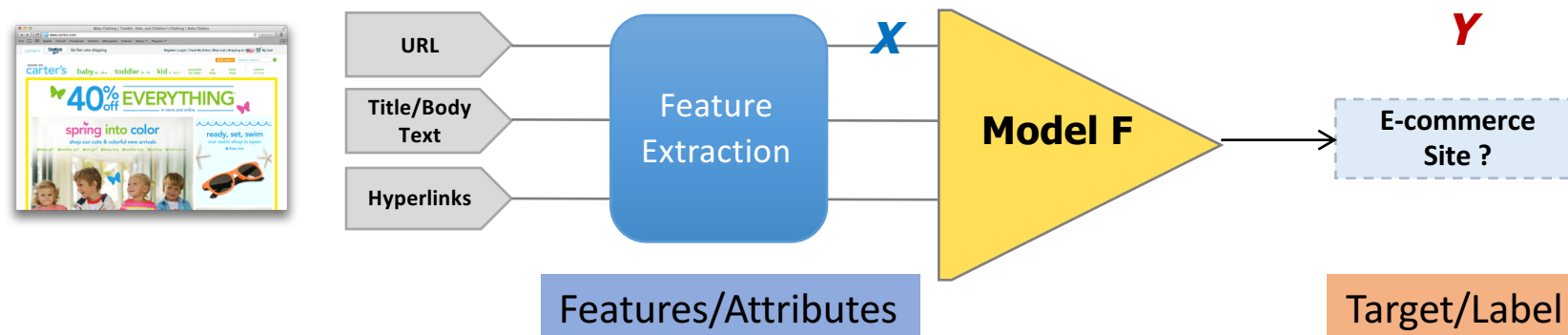


Supervised Learning

- **Training:** Given training examples $\{(X_i, Y_i)\}$ where X_i is the feature vector and Y_i the target variable, learn a function F to best fit the training data (i.e., $Y_i \approx F(X_i)$ for all i)



- **Prediction:** Given a new sample X with unknown Y , predict Y using $F(X)$



- **Inverted Problem:** Given a Y determine relevant X using $F^{-1}(Y)$



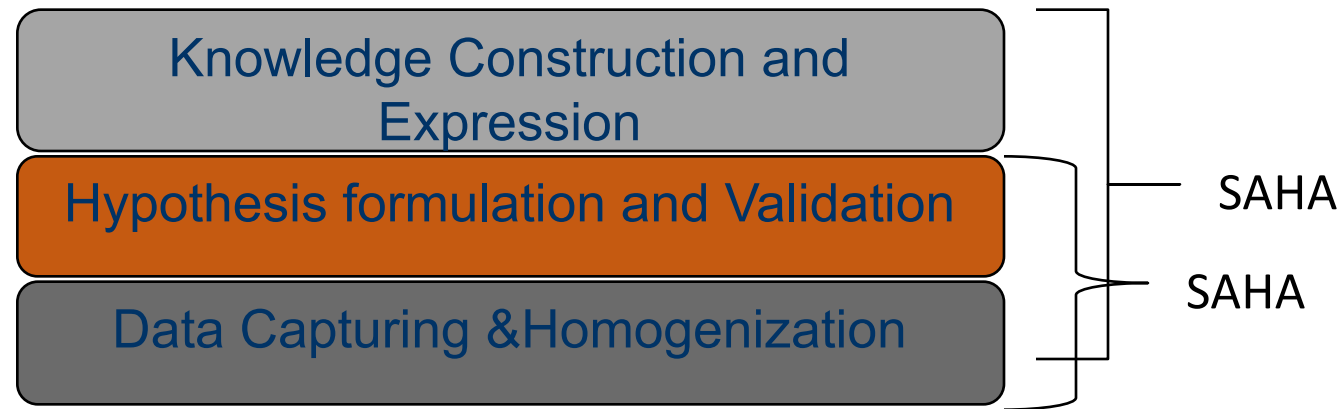
DEXL LAB
EXTREME DATA LAB

SAHA – Sistema de Apoio Holistico ao Atleta

SAHA: Apoio à análise integrada de dados de atletas de alto rendimento



DEXL LAB
EXTREME DATA LAB





- Início
- Atleta/Equipe
- Elementos observáveis
- Métricas
- Rotina
- Coleta de dados
- Gráficos**
 - Análise de trajetória básica (atleta x exame)**
 - Média aritmética e desvio padrão (modalidade x exame x estado)
 - Análise de trajetórias ascendentes e descendentes
 - Análise de trajetórias com extremidades maiores ou menores
 - Análise de máximo e mínimos em uma trajetória
 - Busca trajetórias pertencentes a intervalos máximos e mínimos

Gráficos

Atleta:

Elemento observável:

Data Inicial:

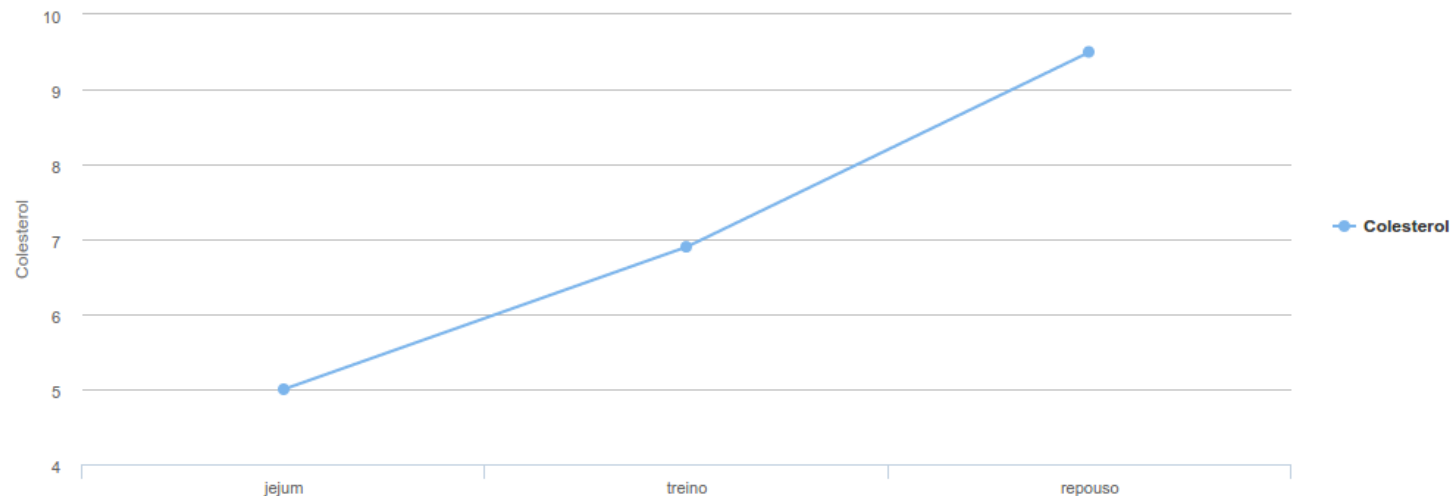
Data Fim:

Estados:

<input checked="" type="checkbox"/> jejum	<input type="checkbox"/> Repouso1	<input checked="" type="checkbox"/> Corrida	<input type="checkbox"/> Corrida 100m
<input type="checkbox"/> Corrida1	<input type="checkbox"/> 12321	<input type="checkbox"/> gdfgfdgdfg	<input type="checkbox"/> 12321
<input checked="" type="checkbox"/> Pos Treino			

Gráfico 1 Gráfico 2 Gráfico 3 Gráfico 4

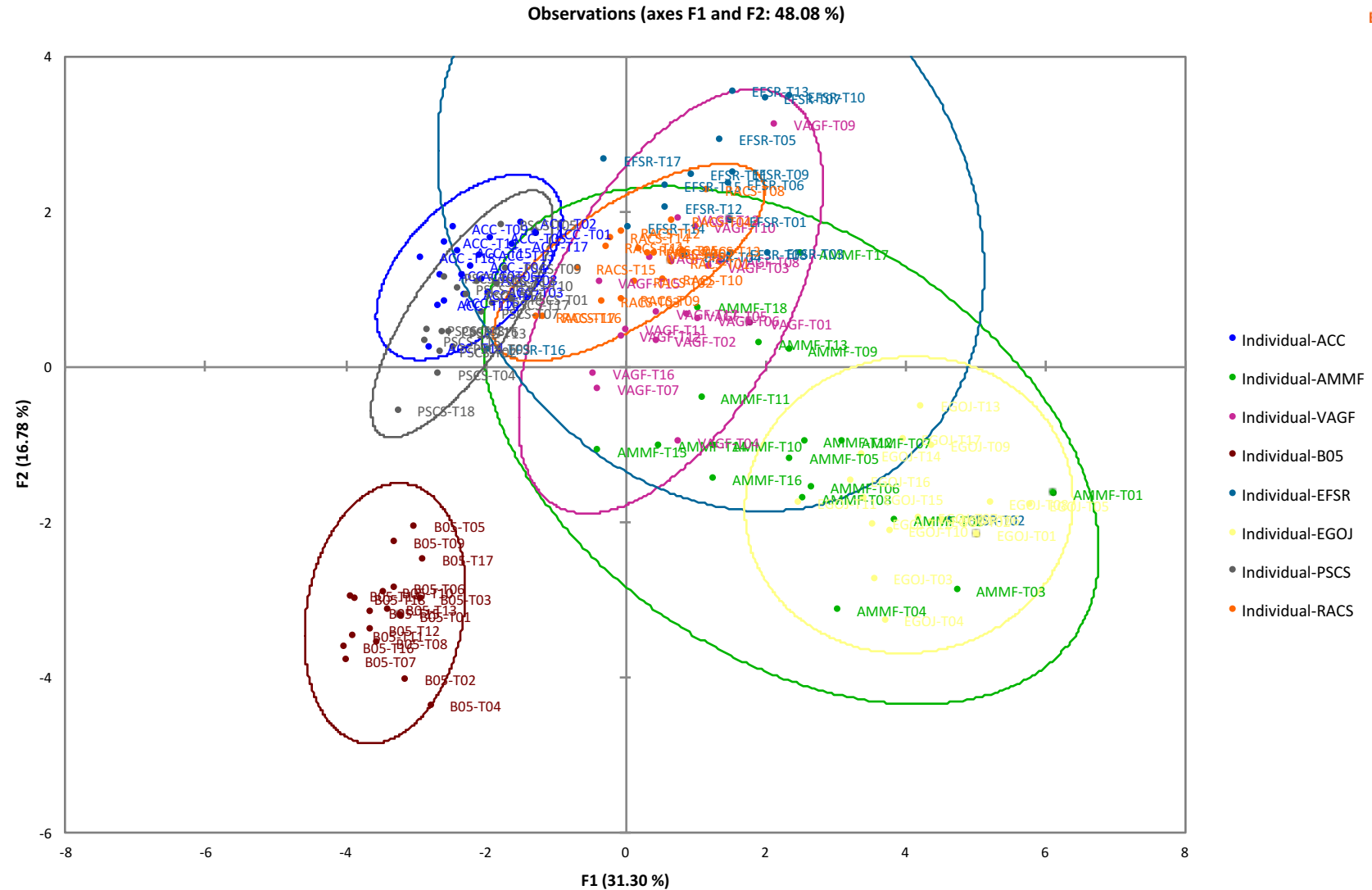
Trajetória básica de: 16/03/2015 a 23/03/2015
(Joao)



PCA: F1 vs F2



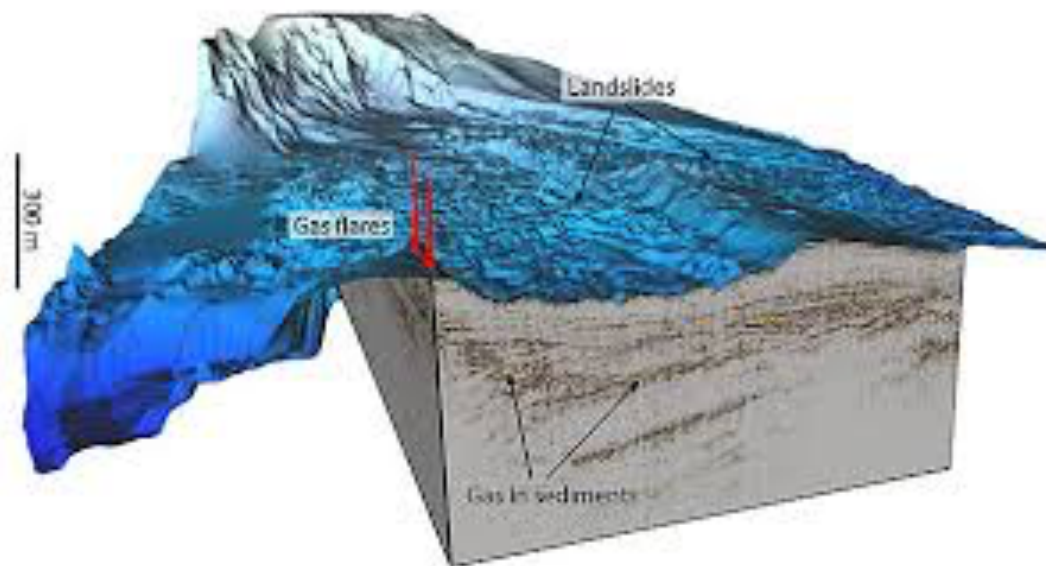
DEXL LAB
EXTREME DATA LAB



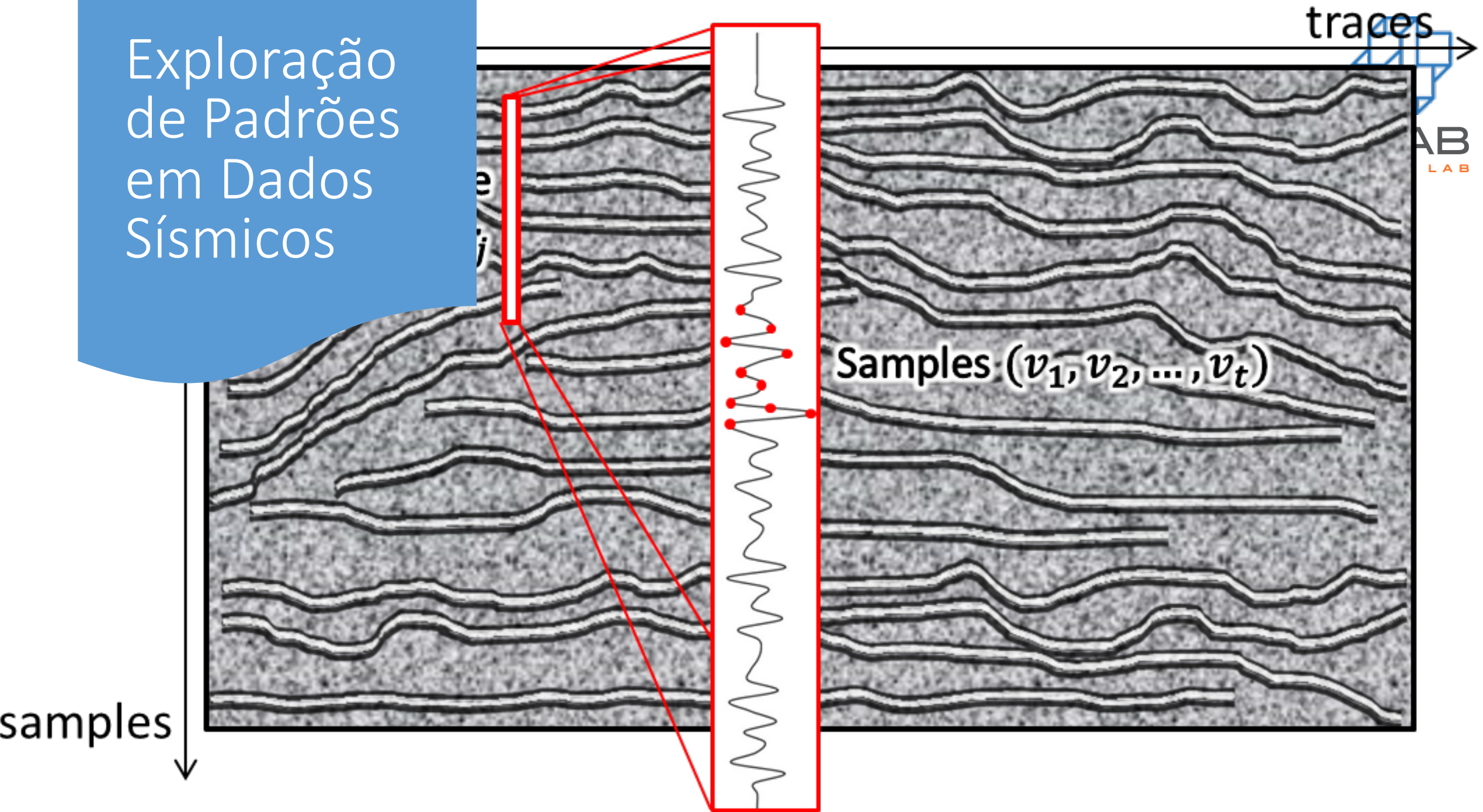
Séries espaço-temporais

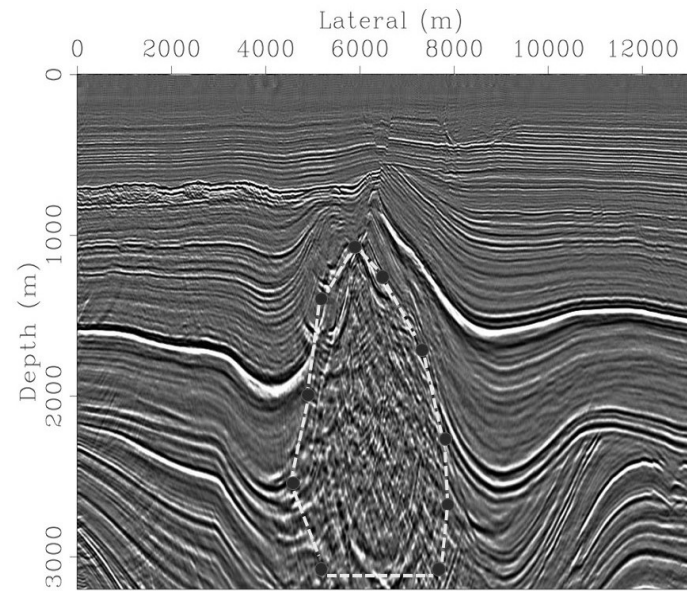
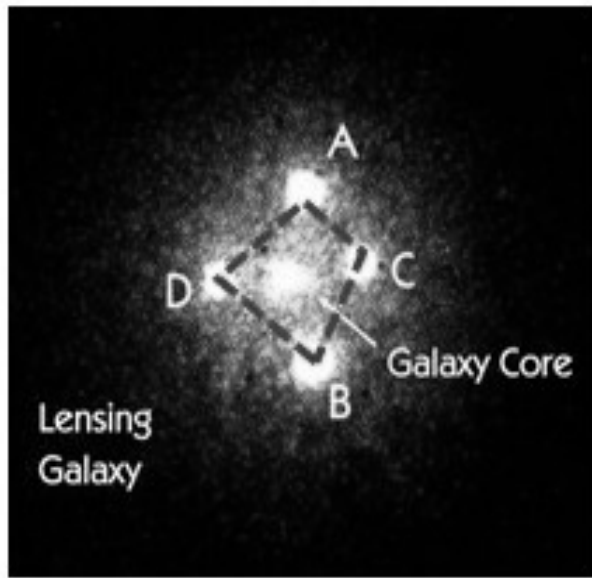


DEXL LAB
EXTREME DATA LAB



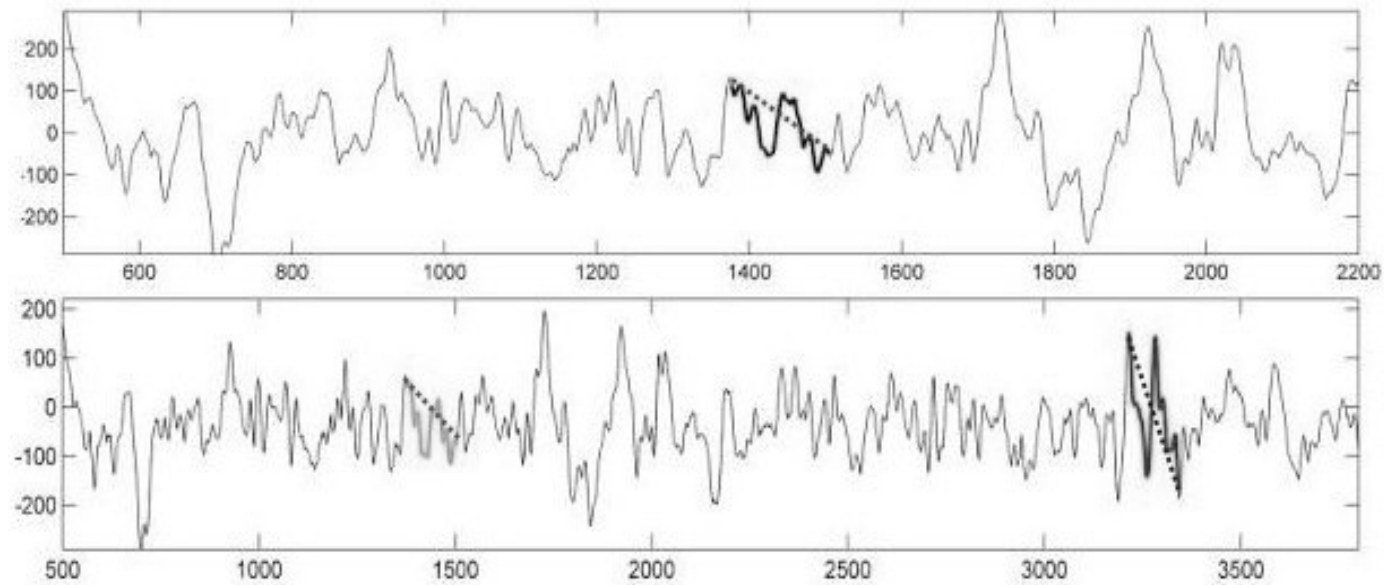
Exploração de Padrões em Dados Sísmicos





(a)

(b)



(c)



DEXLAB
EXTREME DATA LAB

Identificação de padrões em Big Data

- Tratamento dos dados
 - Discretização/normalização
 - Indexação
 - Distribuição
- Análise
 - Algoritmos de análise de padrões em séries temporais
 - Configuração / análise paramétrica
 - Teste e avaliação
- Ranking de soluções
 - Métrica de comparação entre soluções

Comentários Finais



DEXL LAB
EXTREME DATA LAB

- Ciência de Dados
 - Uma nova área multi-disciplinar
 - Ciência da Computação como um dos alicerces na nova ciência
- O processo de investigação carece de sistemas de apoio
 - Ambiente com várias alternativas mas alguns pilares já estão mais sedimentados
- A importância de trabalhar em problemas reais com dados disponíveis
 - Não invente o seu problema!!!
- Os resultados precisam ser validados e interpretados
 - Obtenha e estabeleça um *golden standard*
 - Estabeleça o critério de avaliação para previsões, principalmente sobre dados novos
- Modelos só são bons quando coincidem com a interpretação que se pretende extrair dos dados
- Jornada Em Ciência de Dados, Fev/2016

Um Excelente Início de Curso !!!



Obrigado !😊

Fabio Porto (fporto@Incc.br)

