# Platform Transition: New Technologies and Opportunities

Michael J. Kurtz & the ADS Team mkurtz@cfa.harvard.edu

ADS Users Group Meeting - 11/3/2017







## Change

- 1992 The first Digital Libraries were just being started, conceived, ADS first public demonstration
- 2017 The CFA Library now contains NO BOOKS
- 2042 ?????

### **ADS** Data

- A bibliographic database of 12M scientific publications in Astronomy and Physics (metadata)
- A full-text archive of 4.8M documents (XML, PDF, Latex, OCR text)
- A citation graph with over 6M nodes and 100M edges
- (Anonymous) usage data for 50k regular users

# **Machine Learning, Computational Linguistics**

- Conditional Random Fields
- Latent Dirichlet Allocation (FS; O,T&S 1957, Topic Models)
- Word2Vec
- Various Ad Hoc

## The new way: Machine Learning (CRF)

The AAPS Journal, Vol. 13, No. 2, June 2011 (@ 2011) DOI: 10.208612349-011-9390-2

Research Article

Evaluation of a 2-Integrin Expression as a Biomarker for Tumor Growth Inhibition for the Investigational Integrin Inhibitor E7820 in Predinical and Clinical Studies

### Ron J. Keizer,<sup>14</sup> Y. Funahashi<sup>2</sup> T. Semba<sup>3</sup> J. Wanders<sup>4</sup> J. H. Beijnen,<sup>16</sup> J. H. M. Schellens<sup>45</sup> and A. D. R. Huitens

Received 17 A usuat 2010: accepted 4 February 2011: published online 9 March 2011

Ab shart. E7820 is an ocally active inhibitor of a<sub>cr</sub> integrin mR NA expression, currently tested in phases I and II. We aimed to evaluate what levels of inhibition of integrin expression are needed to achieve turnor stasis in mice and to compare this to the level of inhibition achieved in humans. Tumo con with inhibition was measured in mice bearing a pancreatic KP-1 tumor, dosed at 125-200 mg/kg over 21 days. In the phase I study, 67820 was administered daily for 28 days over a range of 0–200 mg, followed by a 7-day weshout period. PK-PD models were developed in NONMEM. @21ntegrin expression measured on platelets, concesponding to turnor stasis at ⊨21 in 50% and 90% of the mice (In 100, Intern) were calculated. It was evaluated if these levels of inhibition could be achieved in patients at to levable doses. One hundred nineteen 0,2 Integrin measurements and 210 tumor size measurements were available (com mice. The celationship between PK and α<sub>2</sub>-integrin expression was modeled using an indirect-effect model, subsequently linked to an exponential tumo r growth model. Jakes and Jakes were 14.7% (RSE 7%) and 17.9% (RSE 8%). Four hundred sixty two 02 integrin measurements were available form 29 patients. Using the schedule of 100 mg qd (MTD), #2-integrin expression was inhibited more storngly than the J<sub>inteo</sub> and J<sub>inteo</sub> in greater than 85% and greater than 50% of patients, respectively. No decate inhibition of α<sub>2</sub>-integrin expression corresponded to turnor stasis in mice, and similar levels could be reached in patients with the dose level of 100 mg gd.

KEY WORD S bio marker; E7820; modeling and simulation; onco logy; pharmacodynamics.

### INTRODUCTION

Department of Pharmacy and Pharmacology, The Netherlands Cancer Institute/Sio terveart Hospital, Louwesweg 6, 1088 EC, Amsterdam, The Netherlands. <sup>2</sup> Research Labo rate ries Eisai Co., Ltd., Ibaraki, Japan. <sup>6</sup> Eisai Co., Ltd. Hatteld, Hertfordshire, UK. \*Division of Clinical Pharmacology, Department of Medical Oncology The Netherlands Cancer Institute, Amsterdam, The Netherlands. Division of Drug Toxicology, Section of Biomedical Analysis, Department of Pharmaceutical Sciences, Faculty of Science, Utrecht University, Utrecht, The Netherlands. "To whom correspondence should be addressed. (e-mail: con. keise (@sks.nl)

In mammals, 18  $\alpha$  and 8  $\beta$  subunits have been characterized. The investigational anti-cancer drug E7820 is an orally with waving functions related to cell attachment and cell active angiogenesis inhibitor, and has shown anti-tumor signaling. It has been shown in preclinical experiments that activity in several turor models in mice, the effect of which suppression of integrin  $a_2$  by E7830 played a crucial role in was mediated through the inhibition of expression of  $a_2$  inhibition of endothelium table formation (1,2). E7830 was integrin (1,2). Integrins are receptors that mediate attachment evaluated in a phase I dose escalation study in patients with between cells, and between cells and the extracellular matrix makimant solid tumors or lymphomes, the dimical results of (3). They also play an important role in cell signaling. Many which have been reported previously (4) and for which a types of integrins have been identified, multiple types may be population PK model was presented before (5). Since a\_expressed on the cell surface simultaneously. Integrins are integrin is expressed both on turnor cells and on platelets, it has been hypothesized that a s-integrin expression on platelets may be an easily evaluable biomarker for tumor mowth inhibition in response to treatment with E 7820 (1.2).

heterodimens, consisting of  $\alpha$  (alpha) and  $\beta$  (beta) subunits.

The inhibition of a<sub>s</sub>-integrin measured on platelets provides a measure of phanyacological target modulation. and therefore would theoretically provide a better predictor of activity than measures of E7820 plasma exposure. ap-Integrin is therefore currently being evaluated in phase I studies as possible biomarker (46). In this article, we present a modeling and simulation analysis that integrates data from preclinical experiments and a phase I clinical trial to evaluate the expected efficacy of the dose levels that were studies in phase I. Similar analyses have been presented earlier for the development of every linus (7) and getitinib (8). The aim of

2010-71200 (100000-02000) SO 11 The A + the (1). This activity is publicly with a personent at the prime circles a 🛛 😕 🕄

```
🧑 aaos
```

The AAPS Journal, Vol. 13, No. 2, June 2011 (© 2011) DOI: 10 X0X17348-011-0760-7

### Research Article

Evaluation of a .- Integrin Expression as a Biomarker for Tumor Growth Inhibition for the Investigational Integrin Inhibitor E7820 in Preclinical and Clinical Studies

Ron J. Keizer,<sup>3,4</sup> Y. Funakashi,<sup>2</sup> T. Semba,<sup>1</sup> J. Wanders,<sup>4</sup> J. H. Beljnen,<sup>1,5</sup> J. H. M. Schellens,43 and A. D. R. Huttena

> Received 17 August 2010; excepted 4 Petrogev 2011; published online 9 March 2011 Abstract. C. Will be an enable and section in this in a research will N.J. Section 10. Consecting in the section abstract "At and 17.7% (\$55 0%), this incided size two g-otherway transportants why available from 2 that the facto and facto is present that 37% and greater that 20% of partners, respectively. Multi-

KEY WORDS: Strender I This manifer and emploits should be demonstrated.

### INTRODUCTION

Department of Pharmaty and Pharmatelogy, The Netherlands in Multilen in response to irrutation with 8,7520 (1,2). Cancer Institute/Sloterwaart Hospital, Louwewerg G. 1966 EC. Amsterdam, The Netherlands,

Research Laboratories Enai Co., Lui, Buraki, Japan. <sup>3</sup>Eisai Co., Ltd. Hatfield, Hertfordshire, UK.

<sup>4</sup>Division of Chrical Pharmacology, Department of Medical Oncology, The Netherlands Cancer Institute, Amoundam, The Netherlands.

Division of Drug Texasiogy, Section of Biomedical Analysis. Department of Pharmaceutical Sciences, Panalty of Science, Unricht University, Utrechi, The Netherlands.

keizer@hiz.nl)

heterodimers, consisting of  $\Omega$  (alpha) and  $\beta$  (beta) subcritis. In mammab, 18 0 and 8 \$ subunits have been characterized. The investigational anti-cancer drug E.7020 is an orally with varying functions related to cell anachment and cell active anglogenesis inhibitor, and has shown anti-tumor signaling. It has been shown in preclinical experiments that activity in several tarnor models in mice, the offect of which suppression of integrin G, by E.7020 played a crucial role in was mediated through the inhibition of expression of Qr- inhibition of endothelium tube formation (1.2), E7020 was integrin (1,3). Integritis are receptors that mediate attachment evaluated in a phase I down exculation study in patients with between cells and between cells and the extracellular matrix, malignam solid tumors or lymphonus, the clinical results of (3). They also play an important role in cell simulate. Many which have been returned mexicualy (4) and for which a types of integring have been identified, multiple types may be nonpulation PK model was presented before (3). Since G.expressed on the cell surface simultaneously. Integrins are integrin is expressed both on tumor cells and on platelets, it has been hypothesized that G-integrin expression on platelets may be an analy avaluable biomarker for some growth The inhibition of Q2-integrin measured on platelois

provides a measure of sharmainloidcal target modulation. and therefore would theoretically provide a better predictor of activity than measures of E7dJU planma exposure, G-Integrin is therefore corrently being evaluated in phase I studies as monable biomarker (4.6). In this article, we present a modeling and simulation analysis that integrates data from proclinical experiments and a phase I clinical total to evaluate the expected efficacy of the dow levels that were studies in To when correspondence should be addressed. (e-mail: ron. phase I. Similar analyses have been presented earlier for the development of eventimus (7) and gelltinits (0). The aim of

2 880S





### Grobid (http://cloud.science-miner.com/grobid/)

- 11 CRF models (2 for patents)
- Full text processing uses 9 models, 55 final labels, 14 intermediary labels
- More than 10 000 labelled examples





Courtesy: P. Lopez



Download Libraries - Documentation - Examples Community - Developers -

MLlib is Apache Spark's scalable machine learning library.

### Ease of Use

Usable in Java, Scala, Python, and R.

MLlib fits into Spark's APIs and interoperates with NumPy in Python (as of Spark 0.9) and R libraries (as of Spark 1.5). You can use any Hadoop data source (e.g. HDFS, HBase, or local files), making it easy to plug into Hadoop workflows.

data = spark.read.format("libsvm")\
 .load("hdfs://...")

model = KMeans(k=10).fit(data)

Calling MLlib in Python

### Performance

High-quality algorithms, 100x faster than MapReduce.

Spark excels at iterative computation, enabling MLlib to run fast. At the same time, we care about algorithmic performance: MLlib contains high-quality algorithms that leverage iteration, and can yield better results than the one-pass approximations sometimes used on MapReduce.



Logistic regression in Hadoop and Spark

# **Graph/Network Analysis**

- Very complex multi-partite network
- Clustering/Community Detection
- Information Flow, Critical Paths



Download Libraries - Documentation - Examples Community - Developers -

GraphX is Apache Spark's API for graphs and graph-parallel computation.

### Flexibility

Seamlessly work with both graphs and collections.

GraphX unifies ETL, exploratory analysis, and iterative graph computation within a single system. You can view the same data as both graphs and collections, transform and join graphs with RDDs efficiently, and write custom iterative graph algorithms using the Pregel API.

### Speed

Comparable performance to the fastest specialized graph processing systems.

GraphX competes on performance with the fastest graph systems while retaining Spark's flexibility, fault tolerance, and ease of use.

graph = Graph(vertices, edges)
messages = spark.textFile("hdfs://...")
graph2 = graph.joinVertices(messages) {
 (id, vertex, msg) => ...
}

Using GraphX in Scala



End-to-end PageRank performance (20 iterations, 3.7B edges)

## **Network Visualizations**

adsbeta       ouick FIELD: Author       Advanced -       Elbuk query       X       (title)*q       Your search returned 54 results	Feedback O ORCI  ract Year Fulltext All Search Terms  'exoplanet' year:2016 AND read_count:[10 TO { 2 2 Q  Sort: Date desc	D - ☞ Learn - ▲ Account - adsbeta OUICK FIELD: Author First Author	Abstract Year Fulltext All Sean	Feedback (	🕑 ORCID + 🔊 Eearn +	≗ Account +
		Advanced - (title:"exoplanet" year	ar:2016 AND read_count:[10 TO 9999	9999]) 🗶 Q		
Currently viewing data for 54 papers. Change to first papers (max is 54). Submit	Filter current search: X clear select an author or group of authors in the visualization below and click     *add to filter* button	Your search returned <b>201</b> results	8	sort: Date desc \$	Export -	Lini Explore -
Size wedges based on: Author Occurrences Paper Citations	Paper Downloads Summary Detail	Currently viewing data for 201 papers.	▼ Filter current search: x	( clear		
View Link Overlay	Author Network Trails	Change to first papers (max is 201). Culomit Size wedges based on: Number of Paper 2 Paper Citations curves transit exoplanet high atmosphere transmission exoplanet spectrum transit	Group 2 Paper Downloads star exoplanet candidate kepler host	Summary Detail Group 2: atmosphe transit, spectrum, e This group consists of 54 paper times. Top papers from this group in A continuum from clear to clou prinordial water depletion: Sing Characterizing Transiting Exopl Thomas P. (20 citations) Repeatability and Accuracy of I with Post-cryogenic Spitzer. In Detection of H <sub-2< bd=""></sub-2<>	Tre, COCIe, transmi xoplanet Tremove s, which have been cited, clude: dy hot-Jupitar exoplanets f, David K. (64 citations) anet Atmospheres with JV Exoplanet Eclipse Depths I palis, James G. (14 citations) and Evidence to TIO/VO n, and Thomapheres with JV exorplanet Eclipse Depths I palis, James G. (14 citations) and Evidence or TIO/VO n, and Thomapheres with JV exorplanet Settlong I th D189733 b Measured breacher Globads In Exo B. (10 citations) at HD 189733 b Measured Clobads In Exo B. (10 citations) at HD 189733 b Measured (12 citations) at HD 189733 b Measured (12 citations) at HD 189733 b Measured (13 citations) at HD 189733 b Measured breacher Globads In Exo B. (10 citations) at Spoplanets with J2 m 1	ssion, roup from filter in total, 249 without IST; Greene, Veasured in an Ultra- ninant en- with High- ona) planet upper-Earth Class Space-



## **Human-Computer Interface**

- Visualizations
- User Interfaces

## **Paper analytics**



### How it's done: paper to electronic record

•••

adsbeta

+ Back to results

Abstract

Co-Reads

Graphics

Metrics

PT EXPORT

in BibTeX

in AASTeX

in EndNote

Citations (1311)

References (80)

THE ASTRONOMICAL JOURNAL, 124:266-293, 2002 July 000. The American Astronomical Society. All rights merved. Printed in U.S.A

> DETAILED STRUCTURAL DECOMPOSITION OF GALAXY IMAGES CHIEN Y. PENG,<sup>2</sup> LUIS C. HO,<sup>3</sup> CHRIS D. IMPEY,<sup>2</sup> AND HANS-WALTER RIX<sup>4</sup> Received 2001 August 10: accepted 2002 March 27

### ABSTRACT

We present a two-dimensional fitting algorithm (GALFIT) designed to extract structural components from galaxy images, with emphasis on closely modeling light profiles of spatially well-resolved, nearby galaxies observed with the Hubble Space Telescope. Our algorithm improves on previous techniques in two areas: by being able to simultaneously fit a galaxy with an arbitrary number of components and with optimization in computation speed, suited for working on large galaxy images. We use two-dimensional models such as the "Nuker" law, the Sérsic (de Vaucouleurs) profile, an exponential disk, and Gaussian or Moffat functions. The azimuthal shapes are generalized ellipses that can fit disky and boxy components. Some potential applications of our program include: standard modeling of global galaxy profiles; extracting bars, stellar disks, double nuclei, and compact nuclear sources; and measuring absolute dust extinction or surface brightness fluctuations after removing the galaxy model. When examined in detail, we find that even simple looking valaxies generally require at least three components to be modeled accurately, rather than the one or two components more often employed. Many galaxies with complex isophotes, ellipticity changes, and position angle twists can be modeled accurately in two dimensions. We illustrate this by way of 11 case studies, which include regular and barred spiral galaxies, highly disky lenticular galaxies, and elliptical galaxies displaying various levels of complexities. A useful extension of this algorithm is to accurately extract nuclear point sources in galaxies. We compare two-dimensional and one-dimensional extraction techniques on simulated images of galaxies having nuclear slopes with different degrees of cuspiness, and we then illustrate the application of the program to several examples of nearby galaxies with weak nuclei.

Key words: galaxies: bulges --- galaxies: fundamental parameters --- galaxies: nuclei -galaxies; structure - techniques; image processing - techniques; photometric

266

### 1. INTRODUCTION

Galaxies span a wide range of morphology and luminosity, and a very useful way to quantify them is to fit their light distribution with narametric functions. The de Vaucouleurs R1/4 and exponential disk functions became standard functions to use after de Vaucouleurs (1948) showed many elliptical galaxies to have  $R^{1/4}$  light distributions, while Freeman (1970) found later-type galaxies to be well described by a de Vaucouleurs bulge plus an exponential disk. Since then the empirical techniques of galaxy fitting and decomposition have led to a number of notable advances in understanding galaxy formation and evolution. These include investigations into the Tully-Fisher relationship (Tully & Fischer 1977), the fundamental plane of spheroids (Faber et al 1987: Dressler et al. 1987: Diorgovski & Davis 1987: Bender, Burstein & Faber 1992), the mornhological transformation of galaxies in cluster environments (e.g., Dressler 1980; van Dokkum & Franx 2001), the bimodality of galaxy nuclear cusps (Lauer et al. 1995; Faber et al. 1997) and its implications for the formation of massive black holes (Ravindranath, Ho, & Filippenko 2002), and the cosmic evolu-

under NASA contract NASS-20555. 2 Steward Observatory. University of Arizona, 933 North Cherry Avenue, Tucson, AZ 85721; cyp, cimpey@as.arizona.edu. 3 The Observatories of the Carnegie Institution of Washington, 813 Santa Barbara Steecf, Passdem, CA 91101; https://www.edu.

<sup>4</sup> Max-Planck-Institut für Astronomie, Königstuhl 17, Heidelberg, D-69117, Germany.

tion of galaxy morphology (e.g., Lilly et al. 1998; Marleau & Simand 1998)

There are two general types of galaxy fitting: onedimensional fitting of surface brightness profiles (e.g., Kormendy 1977; Burstein 1979; Boroson 1981; Kent 1985; Baggett, Baggett, & Anderson 1998) and two-dimensional fitting of galaxy images (e.g., Shaw & Gilmore 1989; Byun & Freeman 1995; de Jong 1996; Simard 1998; Wadadekar, Robbason, & Kembhavi 1999; Khosroshahi, Wadadekar, & Kembhavi 2000), each with its own tradeoffs and henefits

In one dimension an important consideration is how to first obtain a radial surface brightness profile from a twodimensional image, for which there is no universally agreed upon procedure. A common practice is to use isophote fitting, which is a powerful technique when performed on well-resolved galaxies because it averages over elliptical annuli to increase the signal-to-noise ratio (S/N) at a given radius. However, as many galaxies have isophote twists and changing ellipticity as a function of radius, the galaxy profile is extracted along a radial arc that is ill defined. An alternative approach is to use a direct one-dimensional slice across an image. Burstein (1979) argues that only cuts along the major axis should be used in bulge-to-disk (B/D) decompositions. Meanwhile, Ferrarese et al. (1994) point out that galaxies with power-law central profiles may have different profiles along the major and minor axis.

Fitting profiles in one dimension is frequently used because it suffices for certain goals and is simple to implement. But many studies now resort to two-dimensional techniques. For B/D decompositions a number of authors (e.g., Byun & Freeman 1995; Wadadekar et al. 1999) have



ð 0 = i ui.adsabs.harvard.edu P Feedback D ORCID - Carn - Sign Up Log In QUICK FIELD: Author First Author Abstract Year Fulltext All Search Terms Advanced - author:"^Peng, C\* year:2002 bibstem:"AJ" × Q FULL TEXT SOURCES Detailed Structural Decomposition of Galaxy Images Publisher PDF Publisher Article Show affiliations arXiv e-print Peng, Chien Y.; Ho, Luis C.; Impey, Chris D.; Rix, Hans-Walter E DATA PRODUCTS We present a two-dimensional fitting algorithm (GALFIT) designed to extract structural NED objects (11) components from galaxy images, with emphasis on closely modeling light profiles of spatially SIMBAD objects (11) Archival Data (2) well-resolved, nearby galaxies observed with the Hubble Space Telescope. Our algorithm improves on previous techniques in two areas: by being able to simultaneously fit a galaxy with an arbitrary number of components and with optimization in computation speed, suited for working on large galaxy images. We use two-dimensional models such as the "Nuker" law, GRAPHICS the Sérsic (de Vaucouleurs) profile, an exponential disk, and Gaussian or Moffat functions. The azimuthal shapes are generalized ellipses that can fit disky and boxy components. Some potential applications of our program include: standard modeling of global galaxy profiles: extracting bars, stellar disks, double nuclei, and compact nuclear sources; and measuring Click to view more absolute dust extinction or surface brightness fluctuations after removing the galaxy model. When examined in detail, we find that even simple looking galaxies generally require at least three components to be modeled accurately, rather than the one or two components more SUGGESTED ARTICLES ഒ often employed. Many galaxies with complex isophotes, ellipticity changes, and position angle Early-Type Galaxies in the Sloan Digital twists can be modeled accurately in two dimensions. We illustrate this by way of 11 case Sky Survey, I. The Sample (Bernardi,+); studies, which include regular and barred spiral galaxies, highly disky lenticular galaxies, and more elliptical galaxies displaying various levels of complexities. A useful extension of this algorithm is to accurately extract nuclear point sources in galaxies. We compare two-dimensional and one-dimensional extraction techniques on simulated images of galaxies having nuclear slopes with different degrees of cuspiness, and we then illustrate the application of the program to several examples of nearby galaxies with weak nuclei. Based on observations with the NASA/ESA Hubble Space Telescope, obtained at the Space Telescope Science Institute. which is operated by the Association of Universities for Research in Astronomy (AURA), Inc., under NASA contract NAS 5-26555. The Astronomical Journal, Volume 124, Issue 1, pp. 266-293. Publication Pub Date: July 2002 10.1086/340952 2002AJ....124..266P Bibcode

Keywords Galaxies: Bulges; Galaxies: Fundamental Parameters; Galaxies: Nuclei; Galaxies: Structure: Techniques: Image Processing: Techniques: Photometric; Astrophysics

DOI:

<sup>1</sup> Based on observations with the NASA /ESA Hubble Space Telescape. obtained at the Space Telescope Science Institute, which is operated by the Association of Universities for Research in Astronomy (AURA), Inc. under NASA contract NAS 5-26555.

# Mining, collecting and connecting



### **Future**

- When Computing and Storage are ANOTHER million times more powerful, what will ADS be like?
- ADS is a Human-Machine collaboration
- The New-Gutenbergian Revolution is not yet finished
- ADS is where the state-of-the-art becomes the state-of-the-practice

# **For More Information**

• ADS Bumblebee:

https://ui.adsabs.harvard.edu

• ADS API:

https://github.com/adsabs/adsabs-dev-api

- ADS help and support:
   <u>http://adsabs.github.io/help/</u>
- ADS news:

http://adsabs.github.io/blog/

• ADS users group:

http://adsabs.harvard.edu/adsug.html

adshelp@cfa.harvard.edu

@adsabs