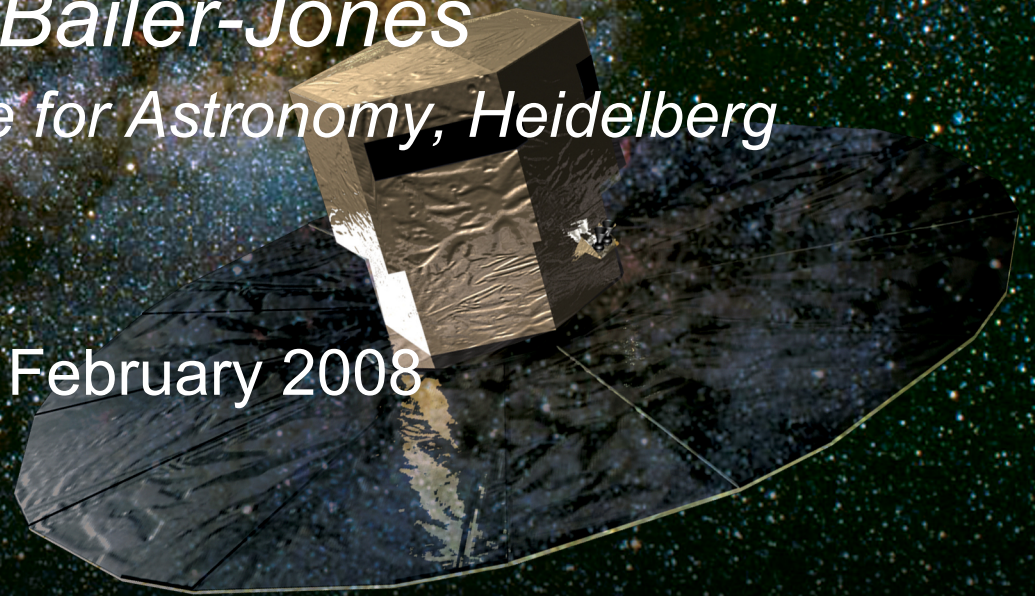# Classification and parameter estimation in astronomy

*Coryn Bailer-Jones*

*Max Planck Institute for Astronomy, Heidelberg*

UCL, 4 February 2008

# Overview

- learning from astronomical data

- Gaia Galactic survey mission

- object classification

- astrophysical parameter estimation

- optimizing the survey: heuristic filter design

- summary

http://www.mpia.de/homes/calj

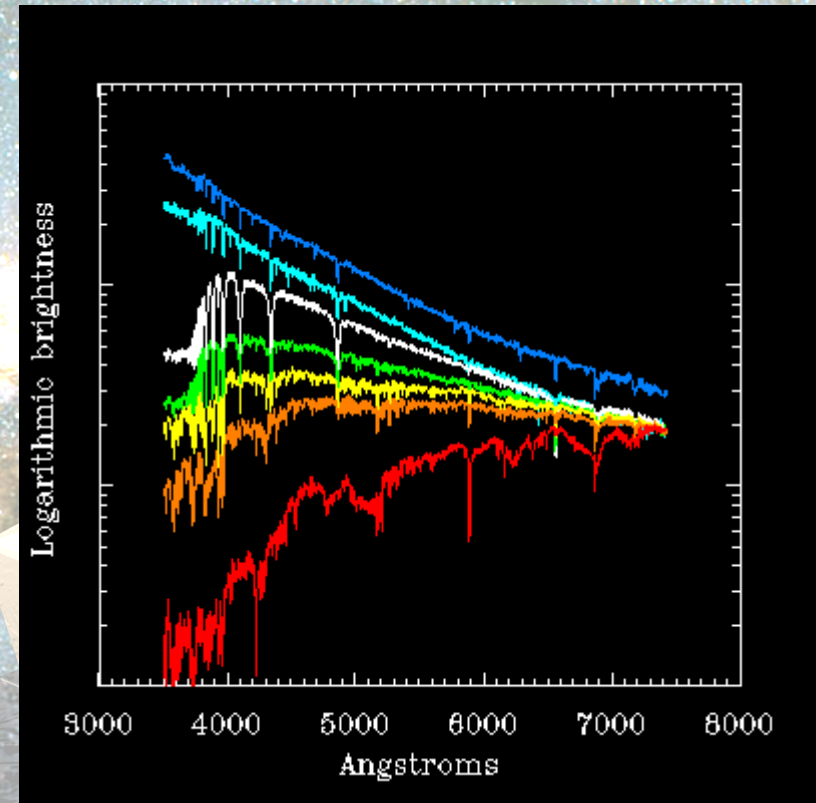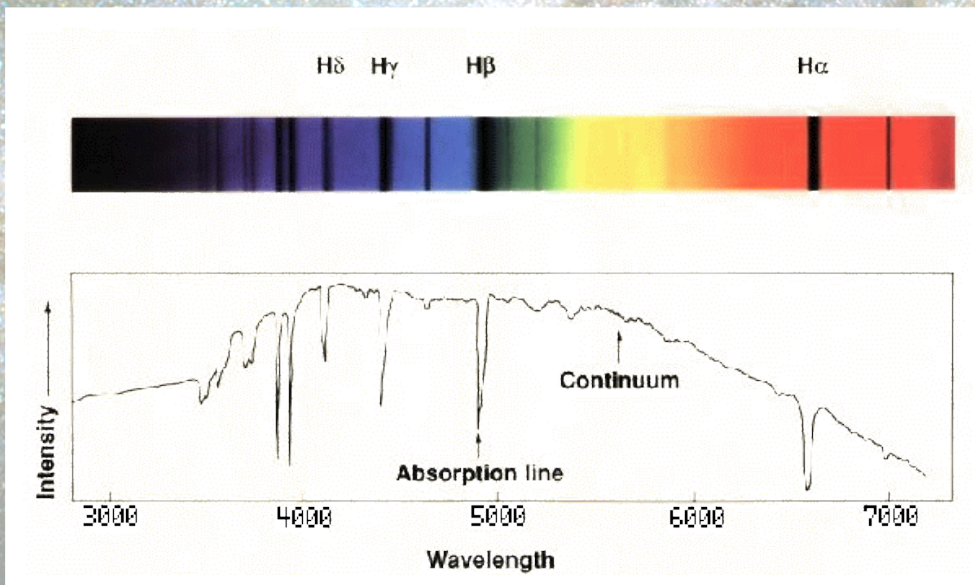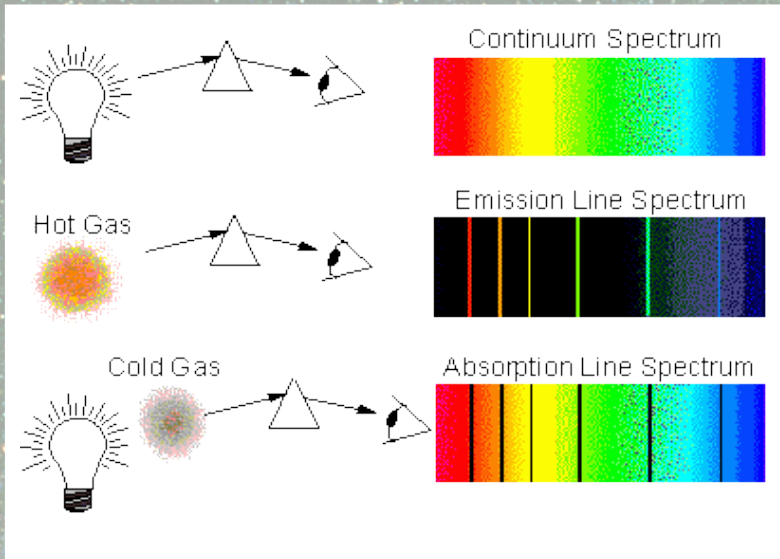Coryn Bailer-Jones, MPIA, Heidelberg

# Acknowledgements

- Gaia @ MPIA group
  - Christian Elting
  - Paola Re Fiorentin
  - Kester Smith
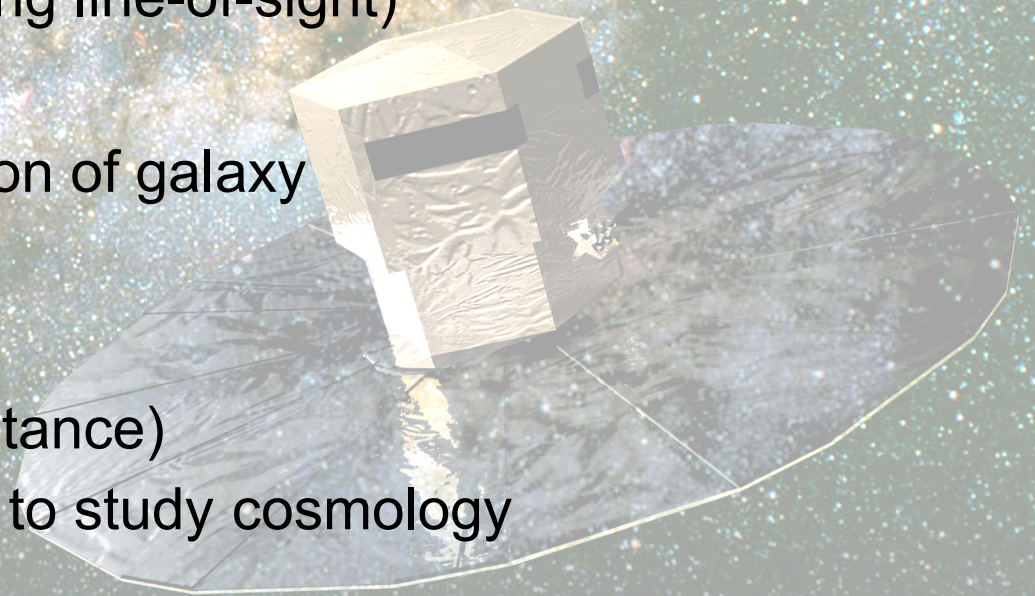  - Carola Tiede
  - Vivi Tsalmantza

- More information on Gaia:
http://sci.esa.int/Gaia

# Astronomical spectra



Coryn Bailer-Jones, MPIA, Heidelberg

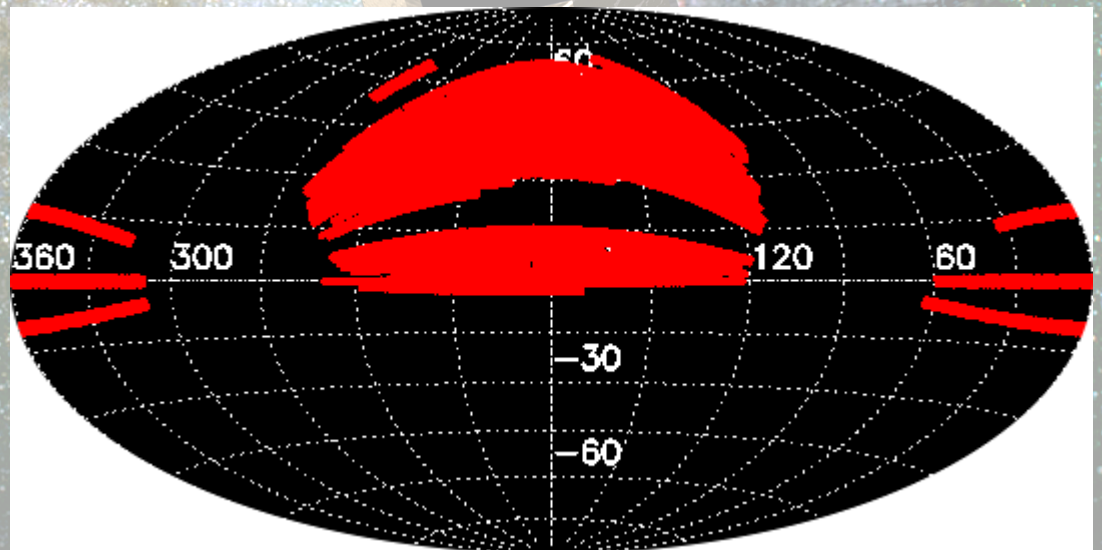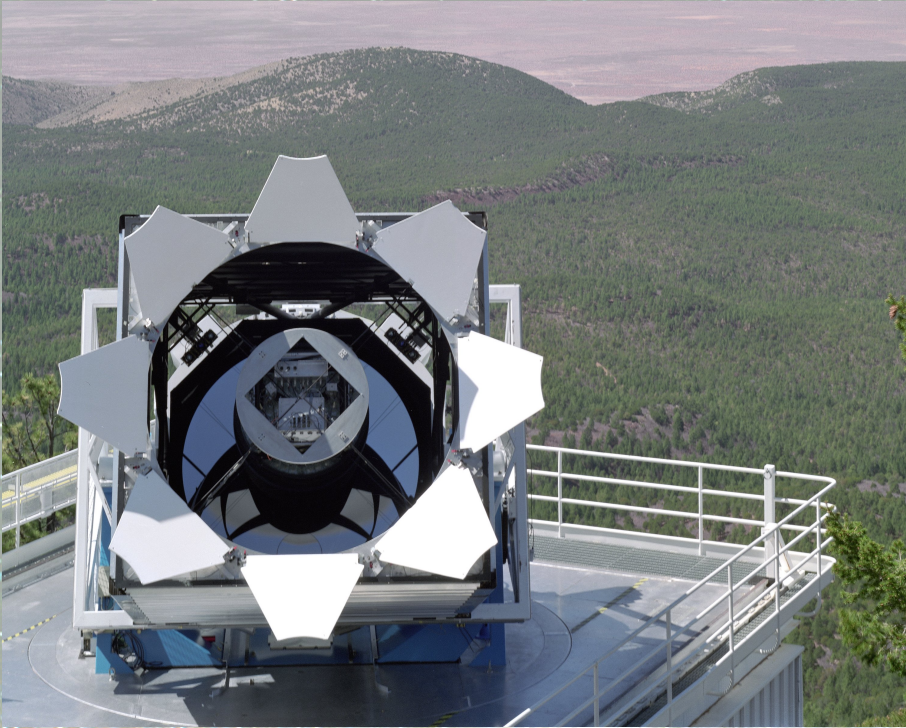images from G. Smith, UCSD and J. Dalcanton, Washington
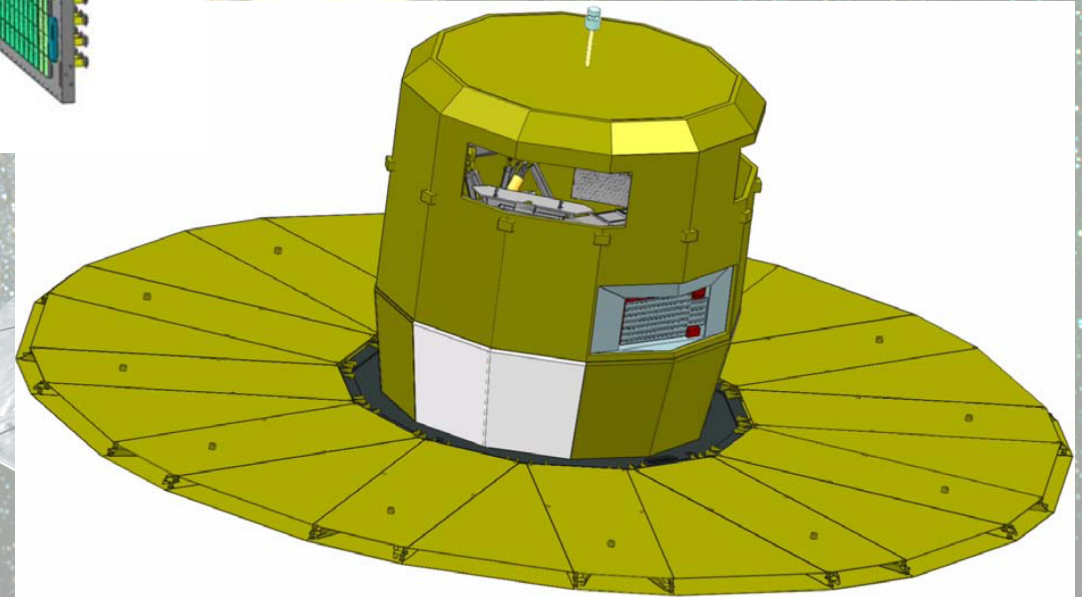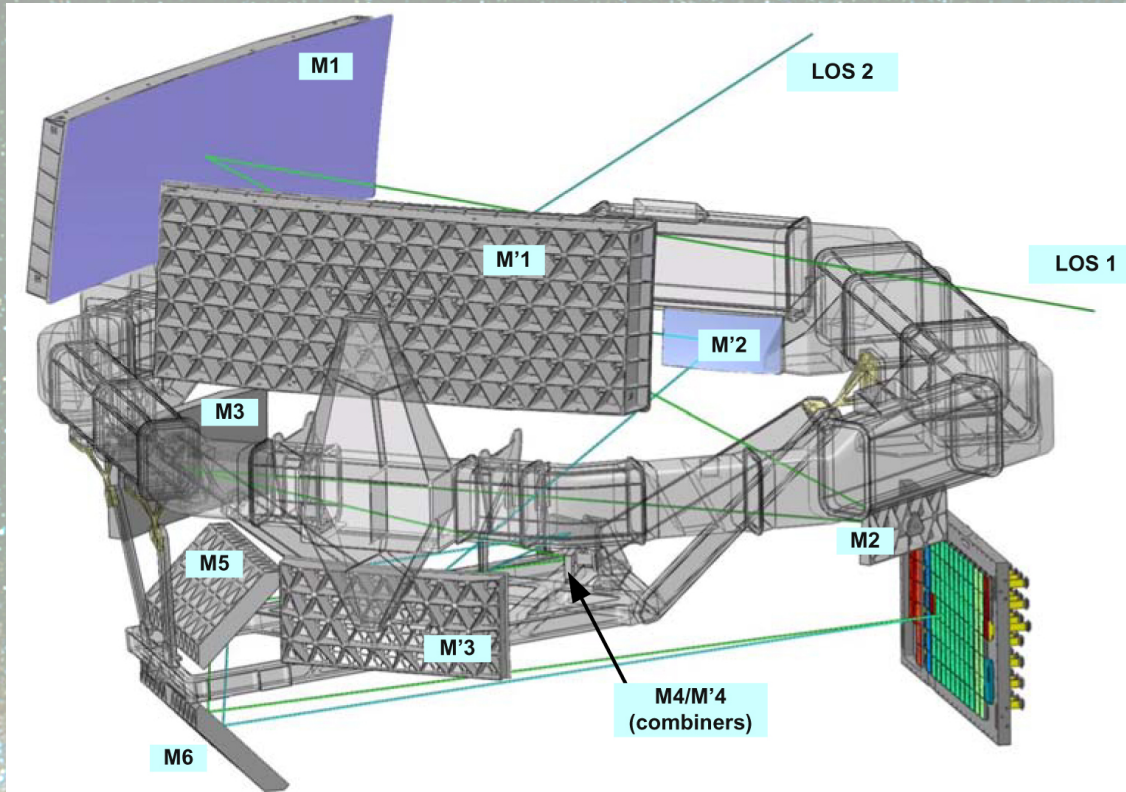
# What can we learn from spectra?

- object type (classification)
  - overall spectral shape
  - characteristic features
- stars
  - atmospheric temperature, pressure, abundances
  - derive mass, age, radius (via evolutionary models)
  - interstellar extinction (along line-of-sight)
- galaxies
  - stellar and gas composition of galaxy
  - infer evolution
- quasars
  - redshift (cosmological distance)
  - use as background lights to study cosmology

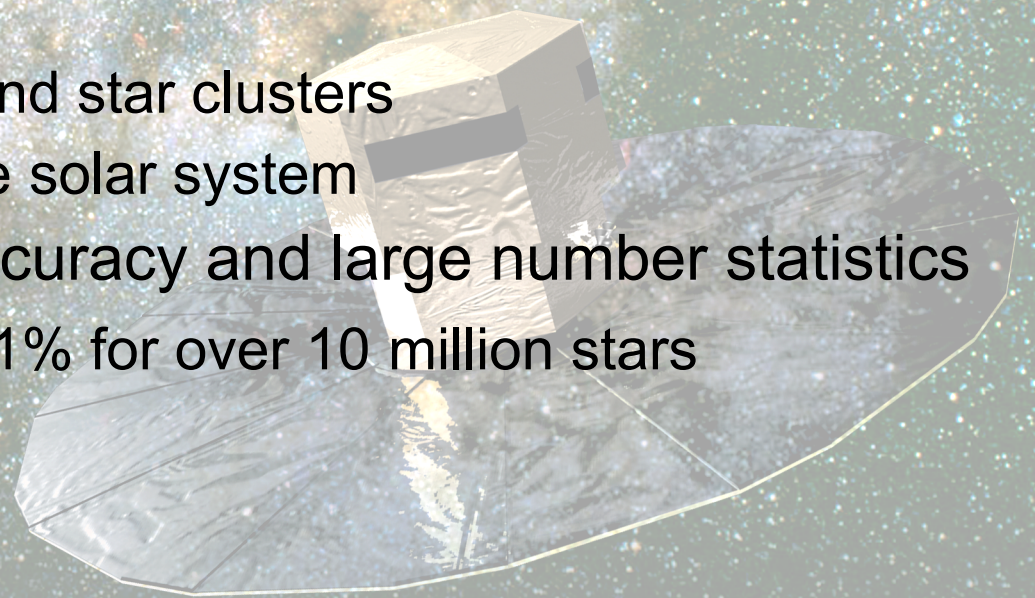Coryn Bailer-Jones, MPIA, Heidelberg

# Astronomical surveys



Coryn Bailer-Jones, MPIA, Heidelberg

# The Gaia Galactic survey mission



M1
LOS 2
M'1
LOS 1
M'2
M3
M2
M5
M'3
M4/M'4
(combiners)
M6

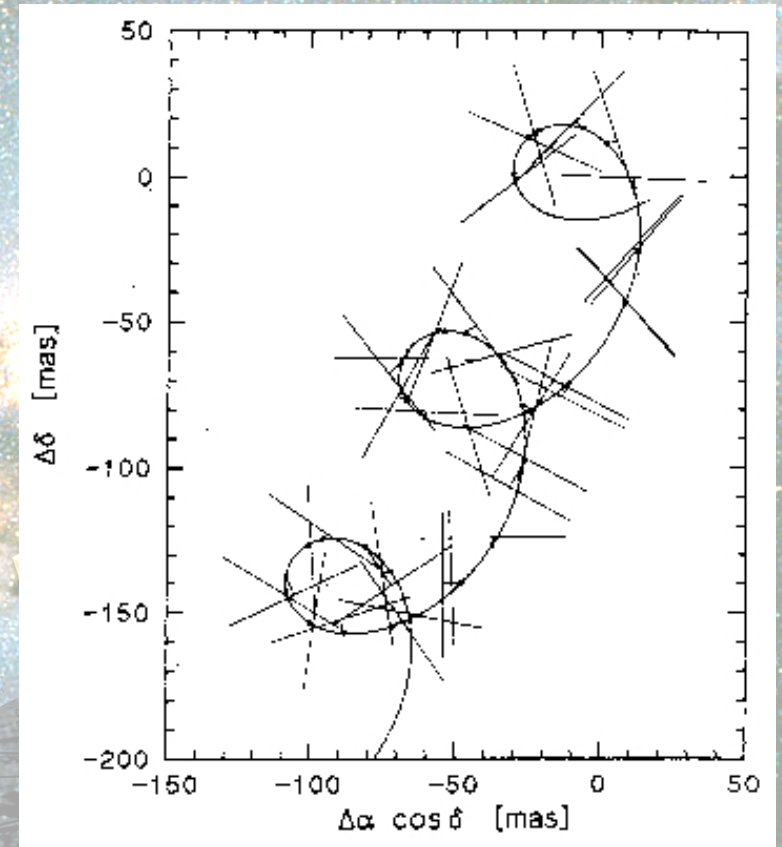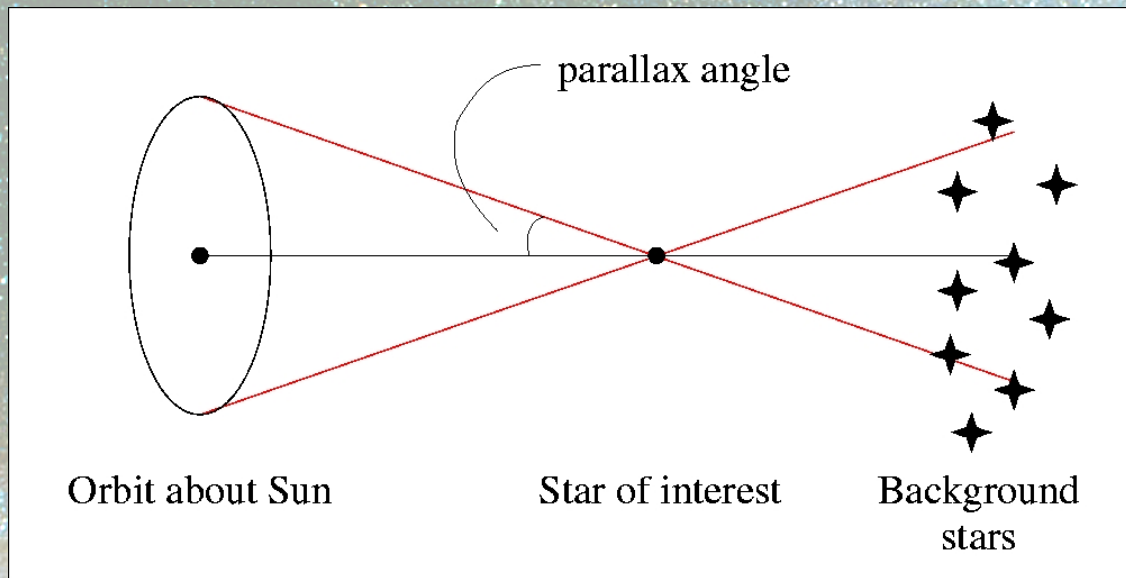Coryn Bailer-Jones, MPIA, Heidelberg

images © EADS Astrium

# The Gaia Galactic survey mission

- five-year all sky survey of one billion objects
- detect objects in real time and measure
  - distances and space velocities
  - spectra (to determine physical properties)
- Science goals
  - structure, origin and evolution of our Galaxy
  - nature of dark matter
  - better understand stars and star clusters
  - find planets outside of the solar system
- strength of Gaia is both accuracy and large number statistics
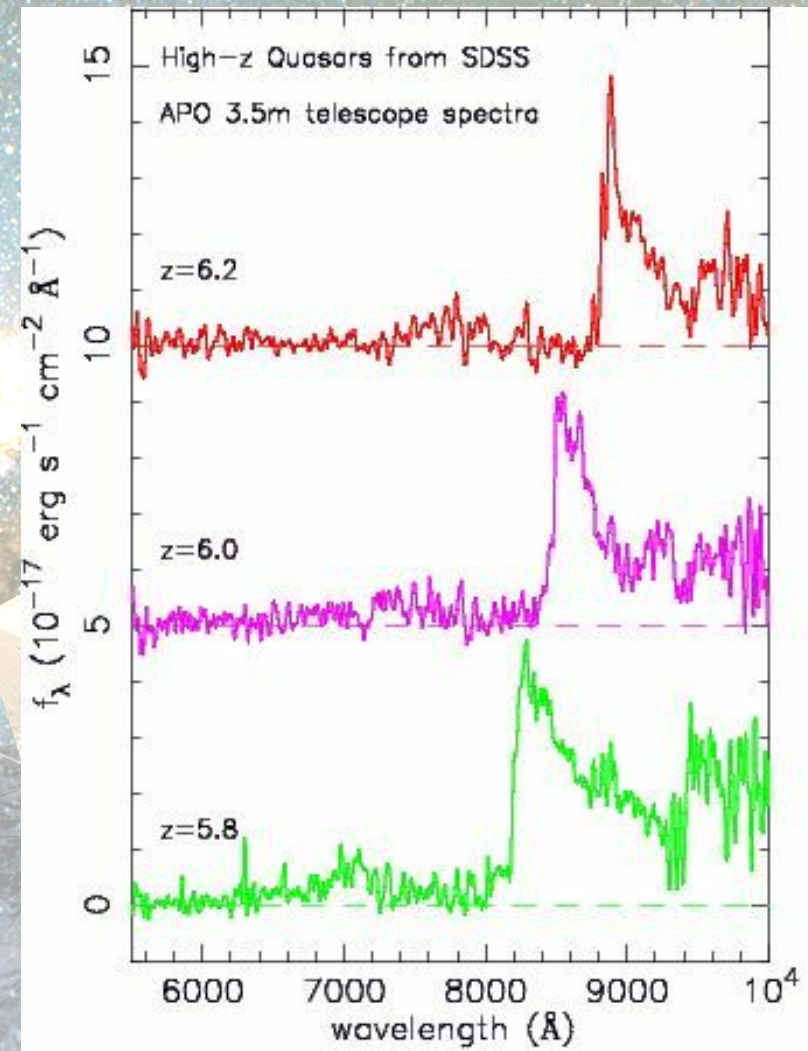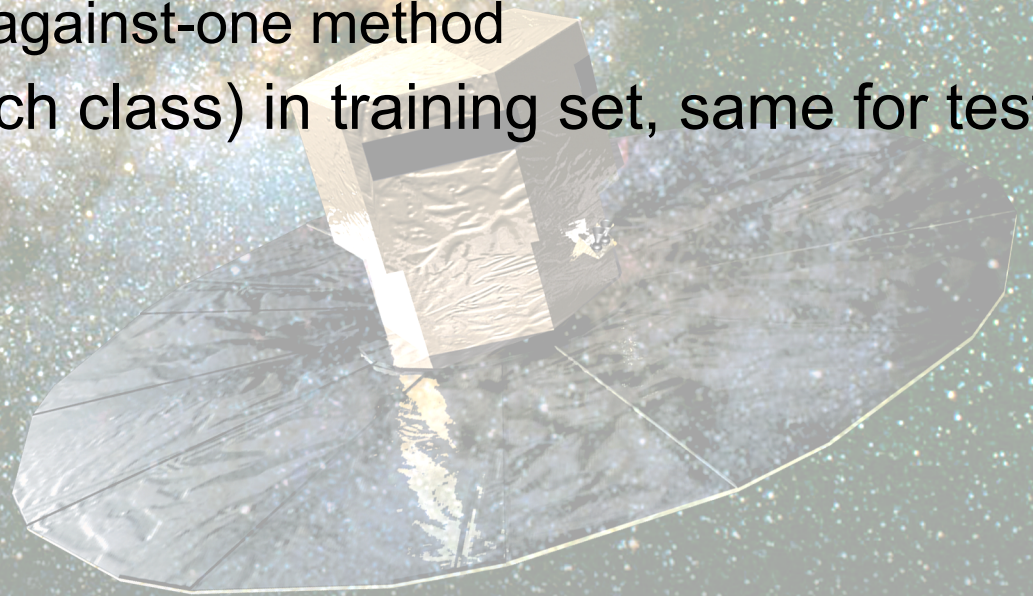  - distances to accuracy of 1% for over 10 million stars

Coryn Bailer-Jones, MPIA, Heidelberg

# Measuring distances via parallax



parallax measurements are relative to reference (background) objects

Coryn Bailer-Jones, MPIA, Heidelberg

# Quasar identification

- extragalactic objects at large distances

- used to define an inertial reference frame ("fixed background")

- characteristic broad emission lines
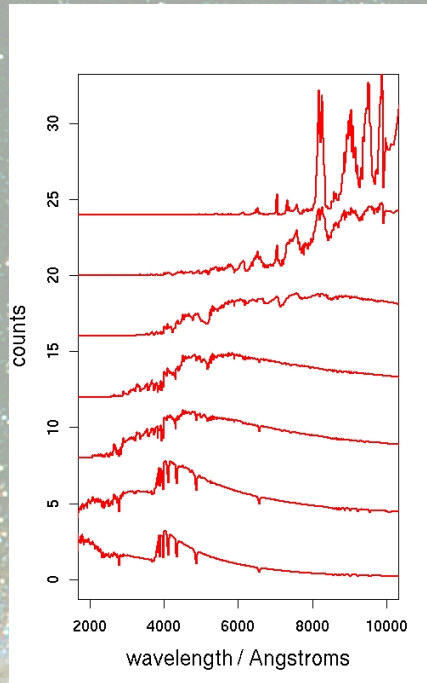
- position of lines in spectrum varies due to redshift



High-z Quasars from SDSS
APO 3.5m telescope spectra

z=6.2
z=6.0
z=5.8

$f_\lambda$ ($10^{-17}$ erg s$^{-1}$ cm$^{-2}$ Å$^{-1}$)

wavelength (Å)

Coryn Bailer-Jones, MPIA, Heidelberg

# The classification problem

- "blind" survey

- object (spectrum) classification via pattern recognition

- four class problem
  - star, quasar, galaxy, physical binary star

- support vector machine classifier (libsvm)
  - RBF kernel
  - multiclass using the one-against-one method

- 8000 objects (~2000 of each class) in training set, same for test
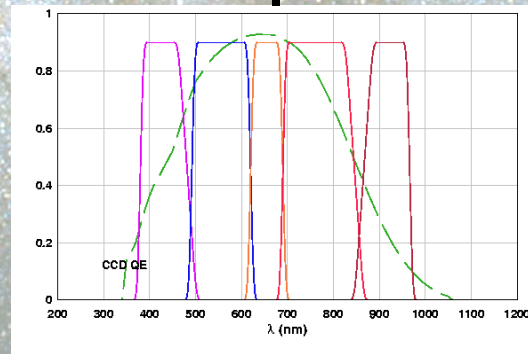
Coryn Bailer-Jones, MPIA, Heidelberg

# Data simulation



synthetic spectra

instrument
simulator
including
noise

instrument model

fluxes
and estimated
covariances

training set

test set

Coryn Bailer-Jones, MPIA, Heidelberg

# Example Gaia spectra



blue: star    red: quasar

Coryn Bailer-Jones, MPIA, Heidelberg

# Gaia spectral data



blue: star   red: quasar      solid: median
dashed: upper and lower quartiles

Coryn Bailer-Jones, MPIA, Heidelberg

# Classification results: contingency table

|          | galaxy | physbin | quasar | star  |
|----------|--------|---------|--------|-------|
| GALAXY   | 99.65  | 0.00    | 0.35   | 0.00  |
| PHYSBIN  | 0.00   | 94.38   | 0.15   | 5.46  |
| QUASAR   | 2.93   | 0.05    | 96.81  | 0.21  |
| STAR     | 0.10   | 13.99   | 3.16   | 82.75 |

- stars and physbin (physical binaries) are hard to distinguish, even in principle

- use output probabilities to vary *completeness* and *contamination* of a sample

Coryn Bailer-Jones, MPIA, Heidelberg

# Classifier confidence

how well does it find the correct objects?

P(class | CLASS)



classifier output probability ⟶

Coryn Bailer-Jones, MPIA, Heidelberg

# Classifier confidence

how well does it find the correct objects?

P(class | CLASS)

how much many false positives in the quasar sample?

P(quasar | CLASS)



classifier output probability ➡️

Coryn Bailer-Jones, MPIA, Heidelberg

# Sample completeness & contamination

select threshold to
build a sample

blue = completeness
red = contamination



threshold probability →

Coryn Bailer-Jones, MPIA, Heidelberg

# Classifier confidence



P(quasar | CLASS)

small but high confidence contamination of quasars by stars

classifier output probability ⟶

Coryn Bailer-Jones, MPIA, Heidelberg

# What is contaminating the quasar sample?



black: stars classified as quasars with high probability

blue: star  red: quasar  solid: median
dashed: upper and lower quartiles

Coryn Bailer-Jones, MPIA, Heidelberg

# Meet reality: population fractions

- in reality quasars are rare compared to stars (ca. 1000 times)

- don't modify training sample distribution/re-train

- instead

  - modify priors
  - treat SVM outputs as likelihoods

Coryn Bailer-Jones, MPIA, Heidelberg

# Meet reality: population fractions

- in reality quasars are rare compared to stars (ca. 1000 times)
- don't modify training sample distribution/re-train
- instead
  - modify priors
  - treat SVM outputs as likelihoods

$$P^{mod}(i) \propto P^{nom}(i) \frac{f_j^{mod}}{f_j^{nom}}$$

$P(i) =$ posterior probabilty that object is in class $i$

$f_j =$ prior fraction of objects in true class $j$

$mod =$ modified, $nom =$ nominal

Coryn Bailer-Jones, MPIA, Heidelberg

# Classifier confidence: nominal vs. modified



Coryn Bailer-Jones, MPIA, Heidelberg

# Modified population fractions (priors)



blue = completeness

red = contamination

relative fractions:

| | |
|---|---|
| star | 1 |
| physbinary | 1 |
| quasar | 0.001 |
| galaxy | 1 |

threshold probability →

Coryn Bailer-Jones, MPIA, Heidelberg

# Modified population fractions (priors)

blue = completeness
red = contamination
solid = nominal
     fractions
     (equal)
dashed = modified
     fractions
     (quasars 1000
     times rarer)



threshold probability →

Coryn Bailer-Jones, MPIA, Heidelberg

# Comparison of different ML models

|          | spectra only |       | spectra + astrometry |       |
|----------|:------------:|:-----:|:--------------------:|:-----:|
|          | K=4          | K=3   | K=4                  | K=3   |
| SVM      | 10.4         | 2.5   | 8.3                  | 0.6   |
| Boosting | 37.8         | 33.9  | 13.0                 | 1.5   |
| MLP      | 9.4          | 1.8   | 7.5                  | 0.2   |
| Mclust   | 10.9         | 0.8   | 9.4                  | 0.1   |
| RBF      | 38.5         | 24.0  | 27.1                 | 15.3  |

overall classification error (percent)
class assignment via highest probability
K=3: 3 classes only (star and physbin merged)

Coryn Bailer-Jones, MPIA, Heidelberg

- several astrophysical parameters (APs) of interest
  - effective temperature, Teff
  - surface gravity, logg
  - abundance, [Fe/H]
  - interstellar extinction, Av
- different nature and degree of impact on the spectra



Coryn Bailer-Jones, MPIA, Heidelberg

Ranges:

$A_v$ = {0, 10}          [Fe/H] = {-2.5, 0.5}

logg = {0.5, 4.5}    $T_{eff}$ = {3350, 35000}

Fixed values:

$A_v$ = 0,          [Fe/H] = 0.0

logg = 4.5,    $T_{eff}$ = 3500

# Strong vs. weak parameters

Minimum distance
(nearest neighbour)
method

**AP ranges:**
$A_v = 0$
$[Fe/H] = 0$
$logg = \{-0.5, 5.5\}$
$T_{eff} = \{2000, 29000\}$



Coryn Bailer-Jones, MPIA, Heidelberg

# An inverse problem



flux in band

log(Teff)

Coryn Bailer-Jones, MPIA, Heidelberg

# An inverse problem



Coryn Bailer-Jones, MPIA, Heidelberg

# Stellar parameters from SDSS spectra

- train regression model using simulated spectra (true physical parameters are known)

- predict (unknown) parameters of real objects

- model: MLP neural network (statnet)

- large input space (thousands of correlated pixels)
  - dimension reduction via PCA
  - also removes some noise
  - can act as a junk filter

Coryn Bailer-Jones, MPIA, Heidelberg

# PCA

original spectra

5+5 reconstruction

25+25 reconstruction

Teff=4752, logg=4.30, [Fe/H]=−1.94

Teff=5251, logg=2.13, [Fe/H]=−0.45

Teff=7720, log g=2.85, [Fe/H]=−0.57

Flux

Wavelength

Re Fiorentin, Bailer-Jones et al. 2007

Coryn Bailer-Jones, MPIA, Heidelberg

# PCA



Coryn Bailer-Jones, MPIA, Heidelberg

Re Fiorentin, Bailer-Jones et al. 2007

Coryn Bailer-Jones, MPIA, Heidelberg

Re Fiorentin, Bailer-Jones et al. 2007

Physical parameter estimates are the basis for detailed physical studies

## ARTICLES

# Two stellar components in the halo of the Milky Way

Daniela Carollo[1,2,3,5], Timothy C. Beers[2,3], Young Sun Lee[2,3], Masashi Chiba[4], John E. Norris[5], Ronald Wilhelm[6], Thirupathi Sivarani[2,3], Brian Marsteller[2,3], Jeffrey A. Munn[7], Coryn A. L. Bailer-Jones[8], Paola Re Fiorentin[8,9] & Donald G. York[10,11]

The halo of the Milky Way provides unique elemental abundance and kinematic information on the first objects to form in the Universe, and this information can be used to tightly constrain models of galaxy formation and evolution. Although the halo was once considered a single component, evidence for its dichotomy has slowly emerged in recent years from inspection of small samples of halo objects. Here we show that the halo is indeed clearly divis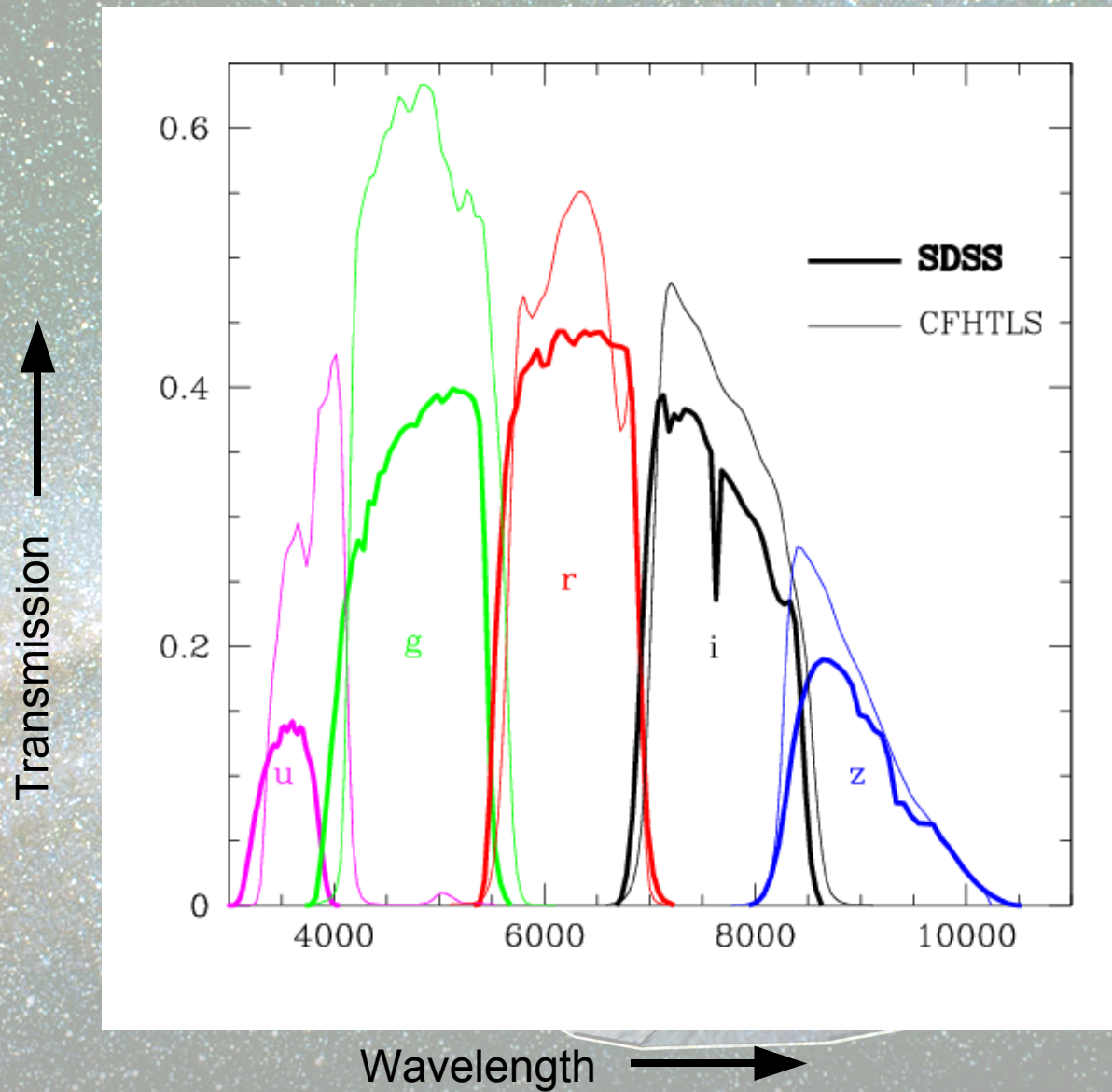ible into two broadly overlapping structural components—an inner and an outer halo—that exhibit different spatial density profiles, stellar orbits and stellar metallicities (abundances of elements heavier than helium). The inner halo has a modest net prograde rotation, whereas the outer halo exhibits a net retrograde rotation and a peak metallicity one-third that of the inner halo. These properties indicate that the individual halo components probably formed in fundamentally different ways, through successive dissipational (inner) and dissipationless (outer) mergers and tidal disruption of proto-Galactic clumps.

Coryn Bailer-Jones, MPIA, Heidelberg

# Incorporating domain knowledge

- much redundant information in spectra

- can directly extract sensitivity of inputs (pixels) from simulated data (physical models)

  - soft dimension reduction

- use in various ways

  - data weighting (global)

  - local regression via local weighting (iterative) *[in progress...]*

- ideal case

  - design optimal filters to use in telescope in first place

Coryn Bailer-Jones, MPIA, Heidelberg

# Photometric filters

image from R. Sharp, IoA

# Heuristic filter design

- algorithm
  - parametrize filter system (three per fiter)
  - establish a figure-of-merit (FoM) of filter system performance
  - maximize FoM with respect to filter system parameters
- optimizer
  - cannot define derivatives
  - local minima guaranteed!
  - evolutionary algorithm (genetic algorithm)
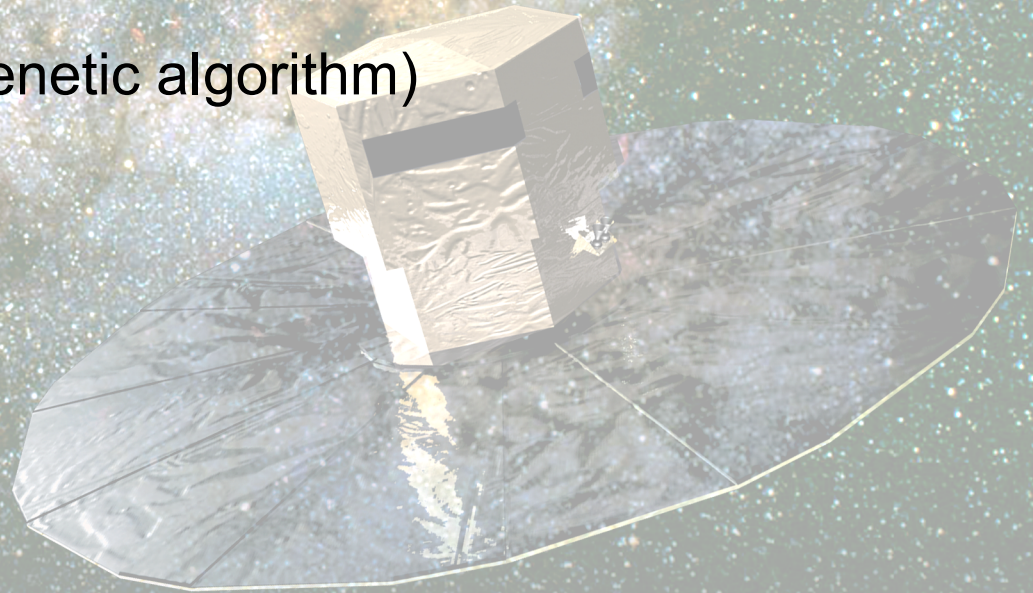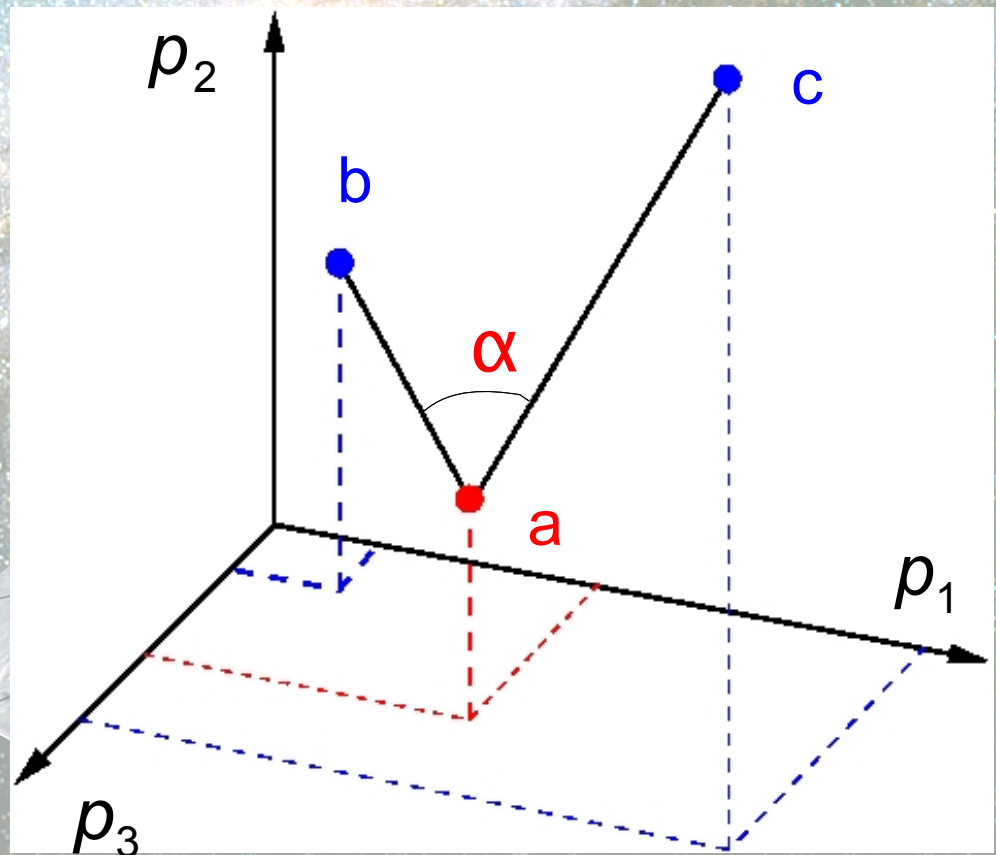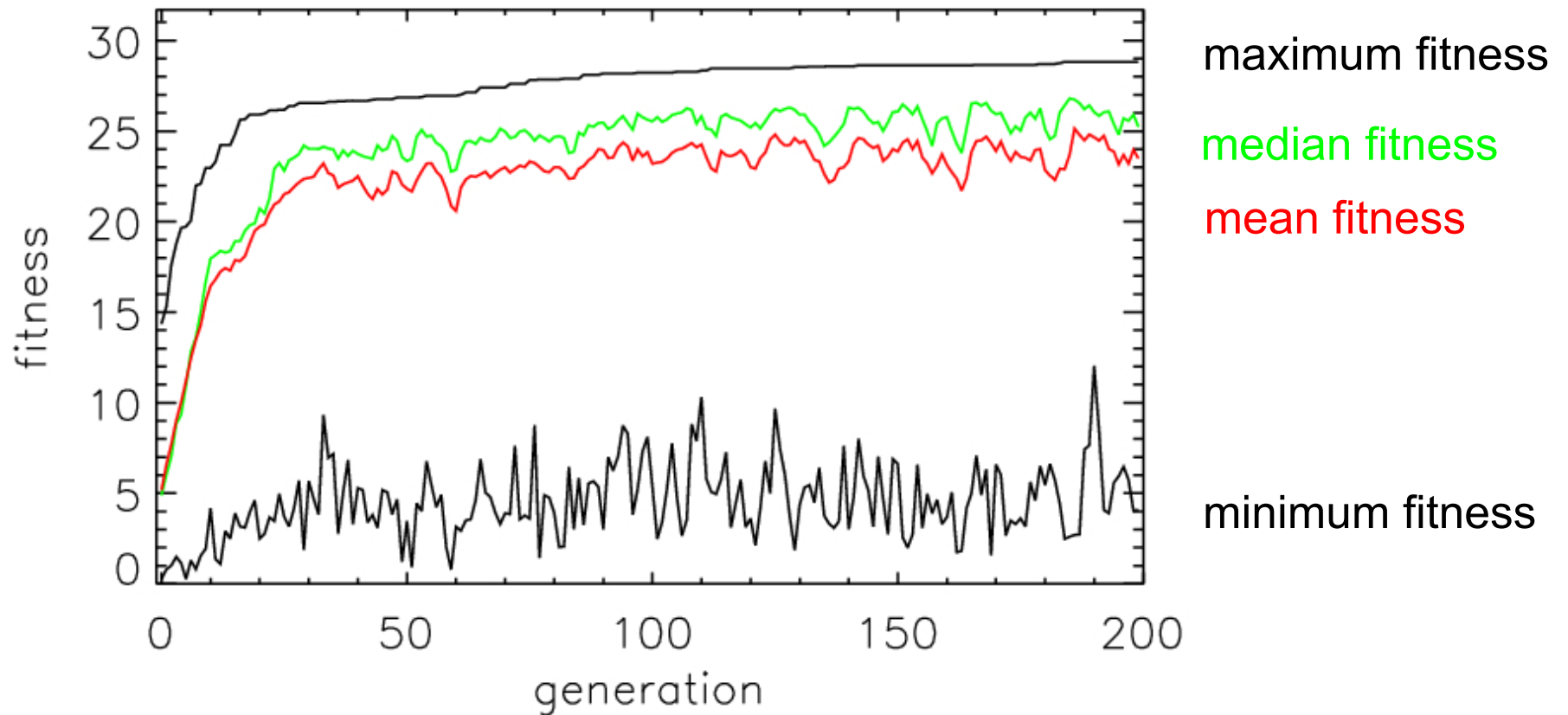
Coryn Bailer-Jones, MPIA, Heidelberg

# Figure of merit (fitness)

- space: flux in each filter of a filter system

- maximum fitness: local directions of variance of the different astrophysical parameters are orthogonal

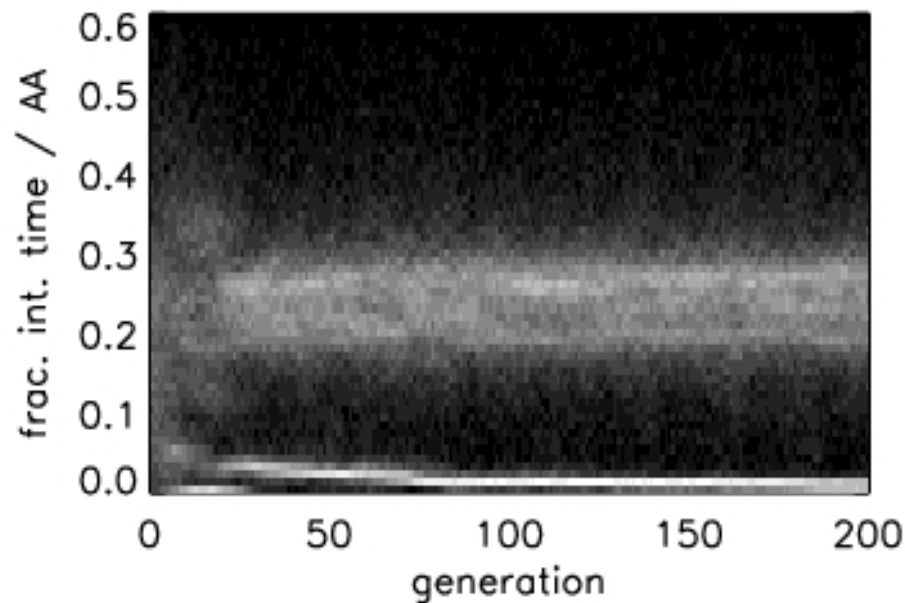- calculated from labelled data set



Coryn Bailer-Jones, MPIA, Heidelberg

# Fitness evolution



maximum fitness
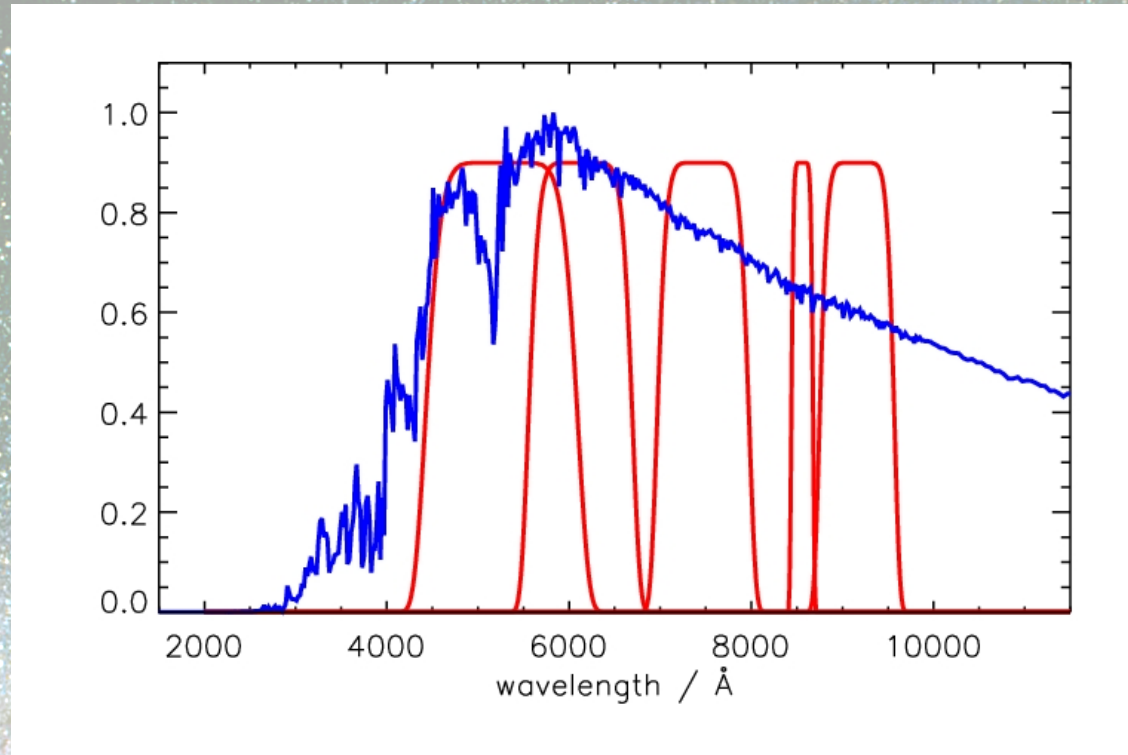
median fitness

mean fitness

minimum fitness

Bailer-Jones 2004

Coryn Bailer-Jones, MPIA, Heidelberg

Evolution of all filter system parameters
(200*5 for each parameter type at each generation)

Bailer-Jones 2004

# Discovered filter system



- **unconventional: overlapping filters**
- **learned filter system competitive in parameter estimation work**
- **method could be extended to design "software filters" in algorithms (e.g. hierarchical classification system)**

Coryn Bailer-Jones, MPIA, Heidelberg

# Summary

- astronomy driven by what we can learn from data (surveys)

- extract optimal information from spectra

  - classification (star, quasar etc.)

  - estimate physical parameters

- use of machine learning methods

  - probabilistic system for classification (SVM)

  - issues of inverse problem and "weak" parameters

- exploit domain knowledge

  - post hoc adjustment based on population fractions (priors)

  - deriving sensitivities from simulated data

  - EA-based heuristic filter design

- future work

  - nonlinear dimensionality reduction

  - degenerate solutions and error estimation

Coryn Bailer-Jones, MPIA, Heidelberg