

# すばるHyper Suprime-Cam (HSC) の 大規模データ処理

高田唯史(国立天文台・天文データセンター)  
HSCデータ解析ソフトウェア開発チーム

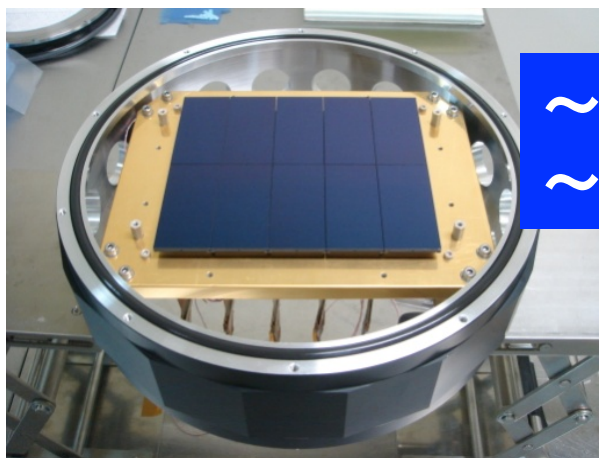
# 話の内容

- 開発の背景(すばるHSCとは、、)
- データフロー(取得->解析->DB)モデル
- 解析手順とアルゴリズム(今回はほぼスルー)
- 解析処理の効率化
- データベースの設計
- 計算機システム(多分スルー)
- まとめ

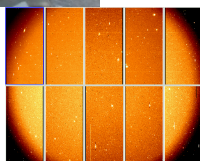
# 開発の背景

# HSC (Hyper Suprime-Cam) とは

- 新しいすばる主焦点可視撮像カメラ
  - 宇宙論(ウイークレンズ)を中心とする戦略的観測
  - 2011年後半にFL予定

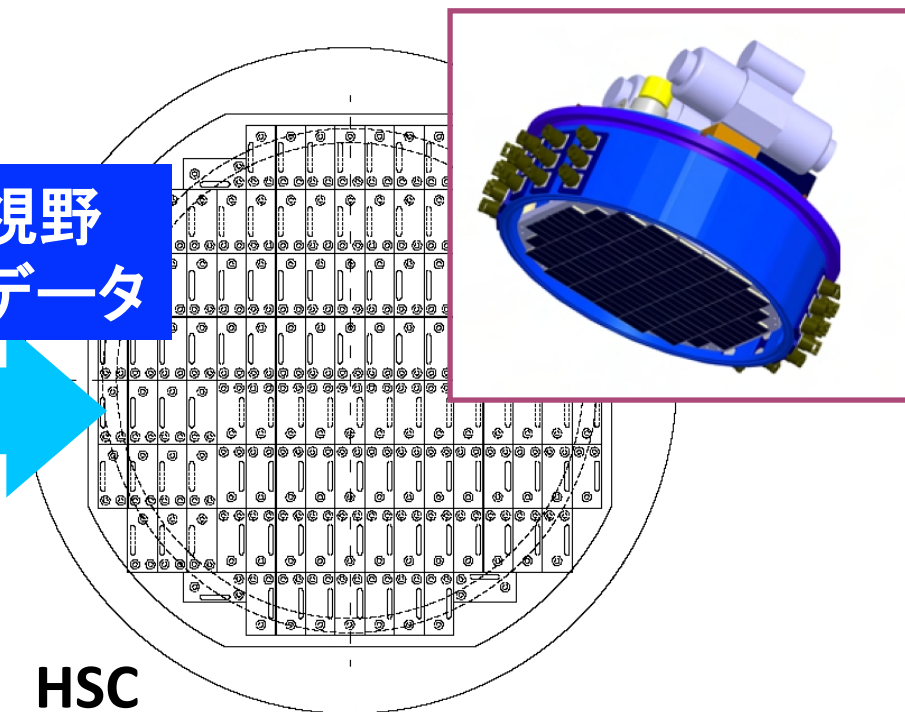


~10倍の視野  
~10倍のデータ



## Suprime-Cam

視野: 34' × 27' (10 × 2k4k CCDs)  
データ量: 185MB/shot (~30GB/夜)  
サーベイ領域: 1~10 平方度

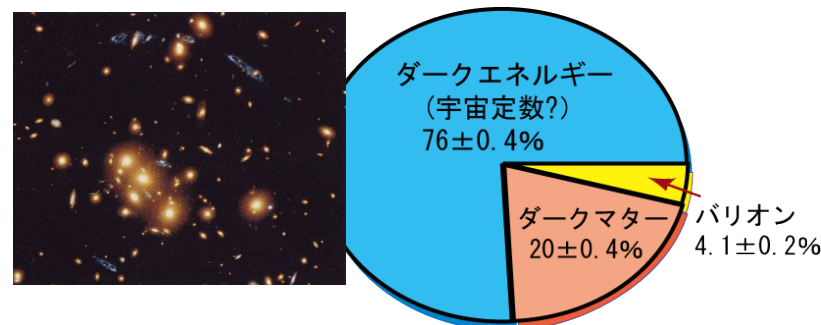


## HSC

視野: 直径1.5度 (104 × 2k4k CCDs)  
2GB/shot (~300GB/夜)  
サーベイ領域 ~2000 平方度

# HSC (Hyper Suprime-Cam) とは

- 特定領域科研費メンバーを中心に、**多岐に渡る興味**を持った研究者の共同研究を予定
  - 宇宙論 (WL、BAO)、超新星、局所・遠方銀河形成進化、突発天体、太陽系天体 etc
- 国内・国際共同開発
  - 国立天文台、  
東京大学/IPMU、KEK、  
Princeton、ASIAA
- (我々が考える)HSCに期待されていること
  - (1) **Suprime-Camと同質かそれを上回る上質な可視Imagingデータを、**
  - (2) Suprime-Camでは到達しえない **広さ&深さ** で得る (e.g., ~1000 sq.deg wide field, >>10 sqdeg deep field)
  - (3) そのデータは**サーベイ観測のランドマーク**となるべき (e.g., HDF, SDSS)



# HSCデータ解析システムの役割

- 期待されている生成物はSCamのものとほぼ同じ
  - 整約済みCCD画像 ( $g', r', i', z', Y$ )
  - モザイク・スタック(ショット合成済み)画像
  - マルチバンド天体カタログ
  - 精度: 位置 $\ll 0.1''$ , 明るさ $\ll 0.05\text{mag}$ , 形状をきちんと把握
- ただし、以下を達成しなければならない
  - 公開に耐える・精度を保証できるプロダクトである
  - 共同研究者へ迅速に提供する
    - 個々人で解析するのはほぼ不可能な規模のデータセット
    - ある程度汎用・共同研究者の必要とする情報を網羅していなければならない

(WL研究→形状、銀河形成研究→深撮像の測光、測光の視野、バンドを通した整合性)

# そのために解決しなければならないこと

基本的にはSCam解析の10倍スケールアップだが、、、

1. HSCというこれまでにない広視野データのための解析  
手順・アルゴリズム

2. 解析の効率化

– 大量のデータを迅速に的確に処理しアップデートするため

3. サイエンスデータベース

が必要。

- これらの完備は我々は未体験！

- もちろんこのほかに、フラックスキャリブレーションソースの確保、データの性質に影響の出る装置状態の把握といった、解析だけに閉じない課題はある

# どのような運用でデータを効率よく解析し、成果の獲得を目指すのか？

(特に期待される大サーベイのデータからの成果の効率よい獲得を目指すには？)



高速かつ精度の高い画像処理とそれに伴うデータのクオリティ・コントロールで無駄のない運用を行う。



マウナケア山頂、ハワイ山麓施設、日本での処理の役割分担を考える





日本



国立天文台:三鷹  
東大IPMU:柏

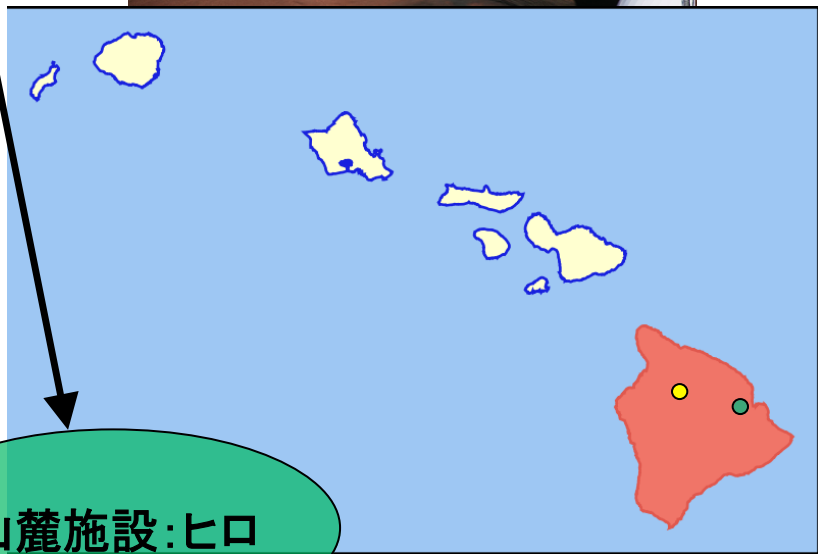
オフライン解析+アーカイブ  
(数日程度のサイクルでフル解析)

マウナケア山頂

観測データ取得

1Gbps

ハワイ

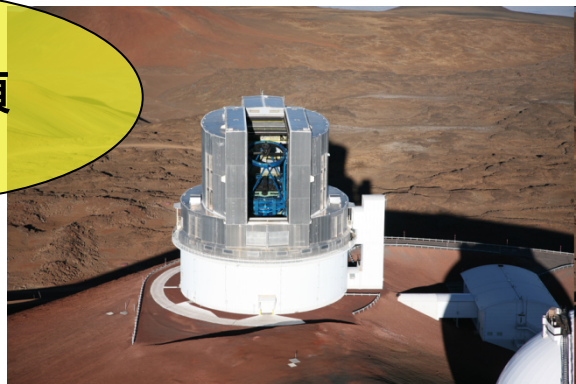


山麓施設:ヒロ

オンライン解析  
(~5分程度のサイクル)

?オフライン解析?  
(数日程度のサイクル)

155Mbps





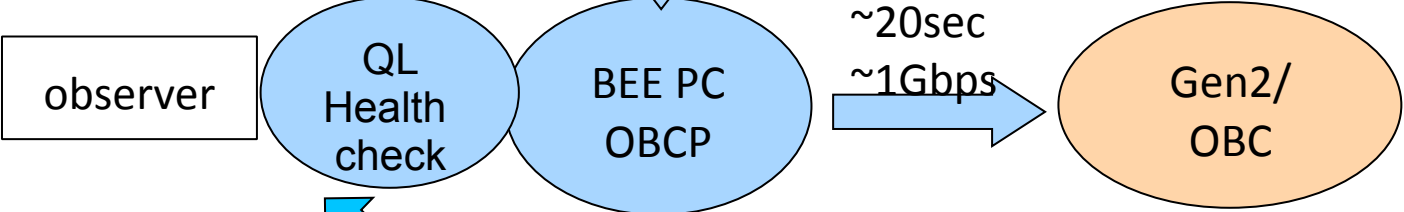
OCS/  
CoDM

HSC

HSC BEE

Data production rate:  
• 2GB/shot  
• ~ 500GB/night  
• ~ 5TB/run  
• ~150TB/300nights

**Summit**



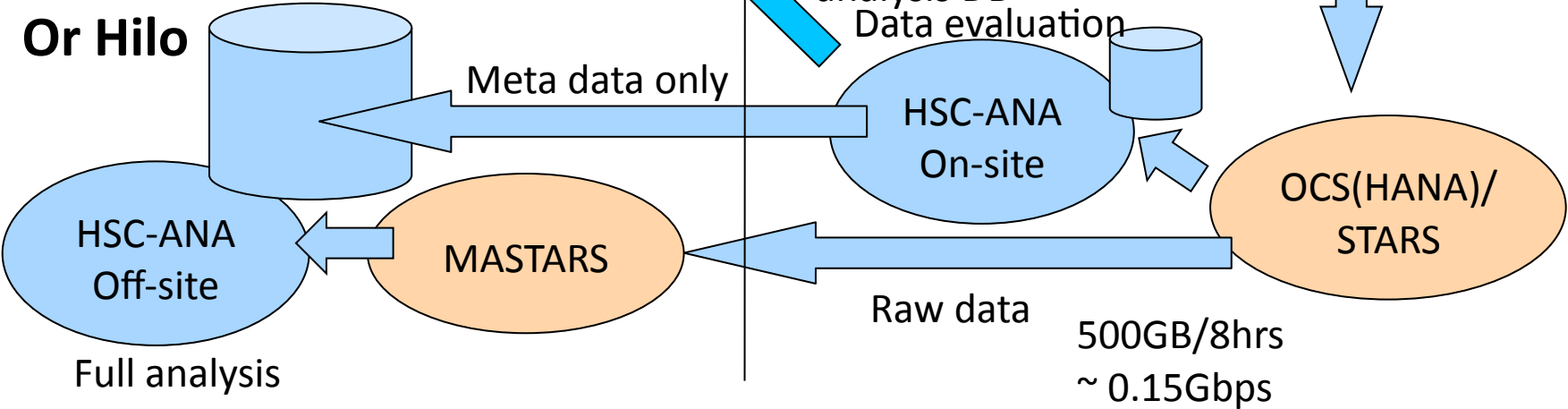
Feedback:  
<< 3min

Summit-Hilo  
1Gbps

**Mitaka  
Or Hilo**

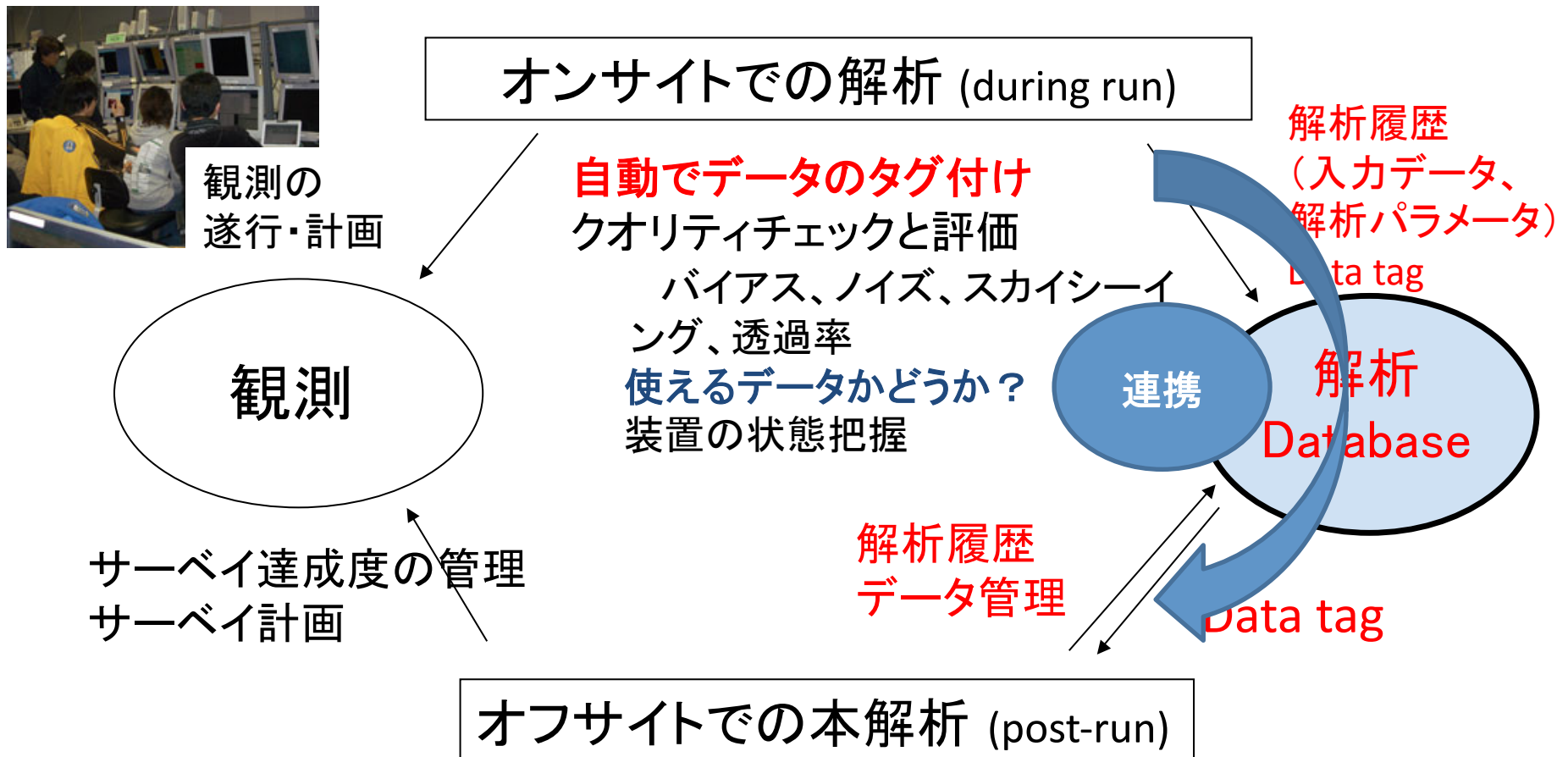
3-run processed data ~30TB?  
+ analysis DB  
Data evaluation

**Hilo**



# データベースを用いた解析管理 オンサイト解析と本解析の連携

- ハワイ観測所にオンサイト解析システムのプロトタイプを構築し、**SCamの観測支援**と、解析の効率化のためのデータベース+ミドルウェアを試験。共同利用観測で運用試験している。

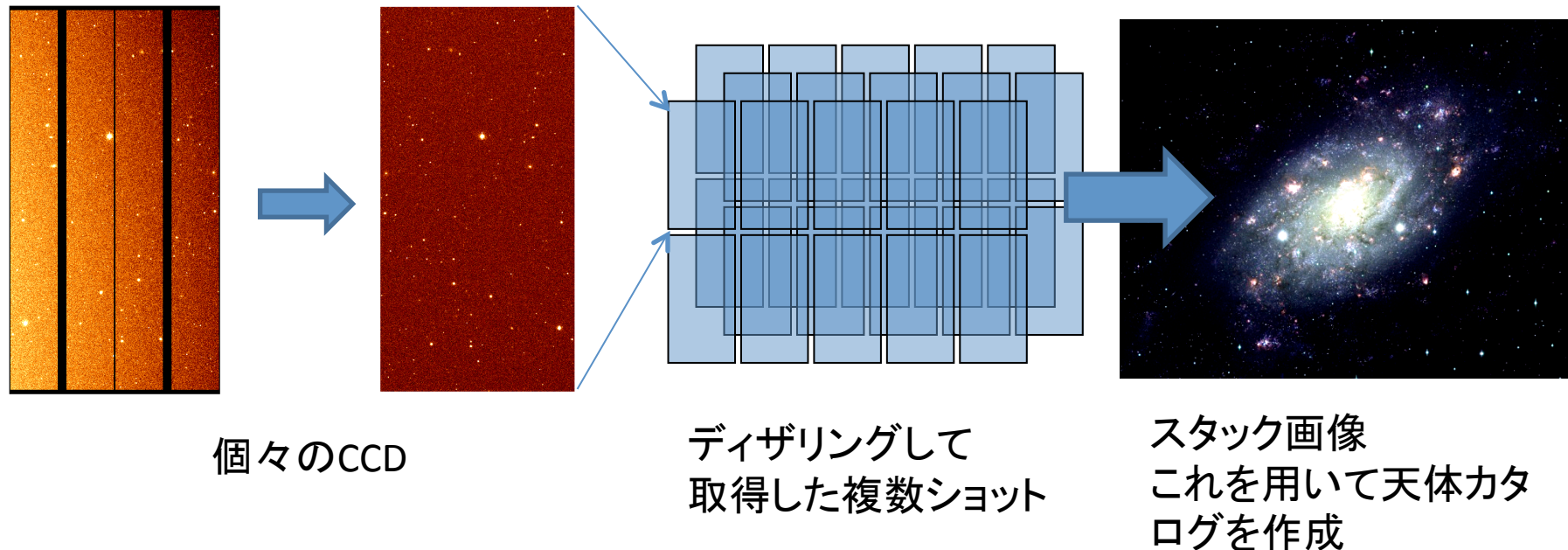


# HSCデータのための解析手順・ アルゴリズムについて

(今回はあまり細かく述べない、、)

# 一般的な撮像データ解析手順＝HSC でもベースとする

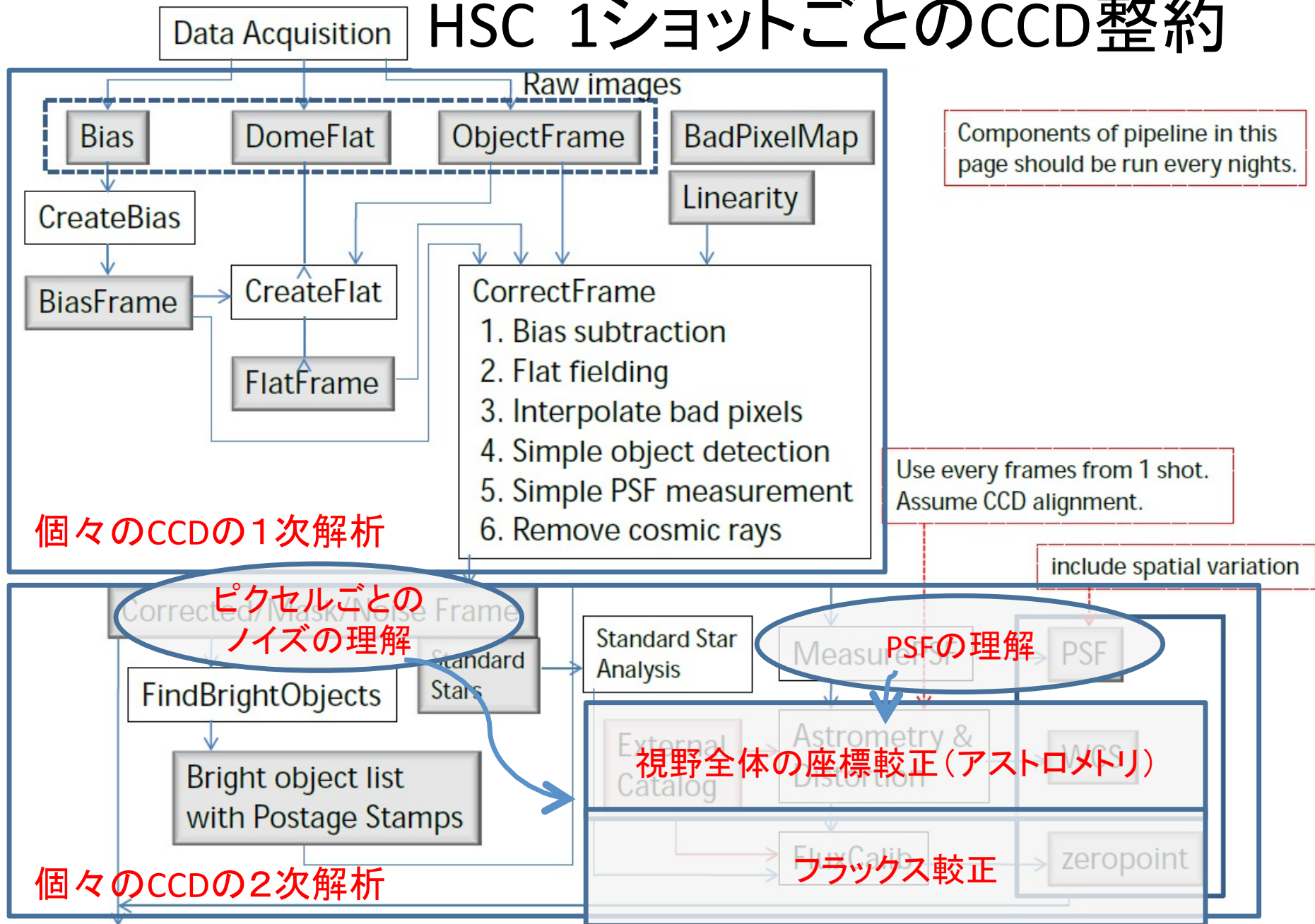
1. 各CCDのバイアス除去・感度差補正（一次処理）
2. 必要な一次処理済みCCDデータを集め、
3. モザイクして足し合わせる
4. 位置情報、カウントを物理量に直す



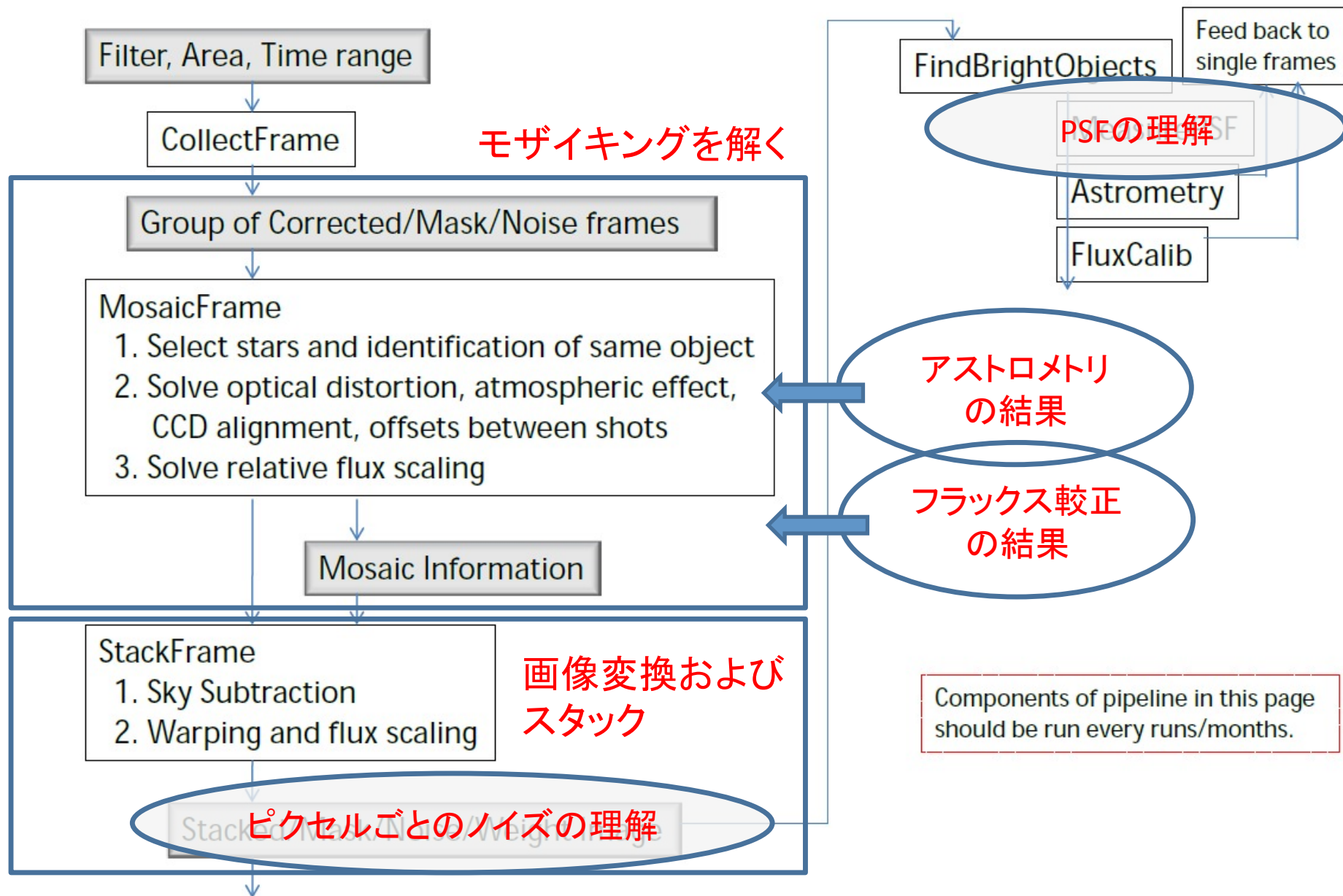
次の3枚のスライドでHSCでターゲット  
とする解析手順を紹介します



# HSC 1ショットごとのCCD整約

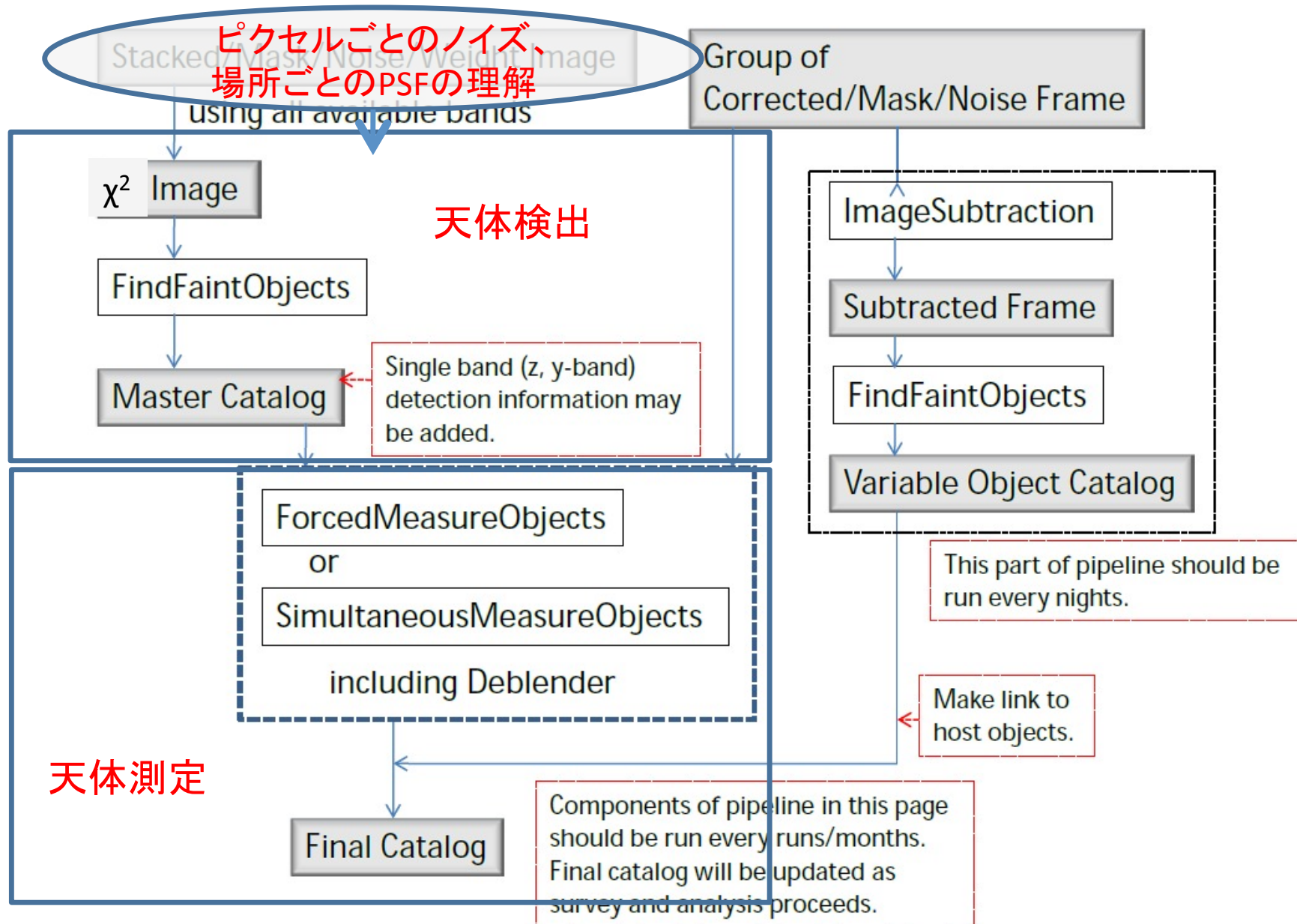


# HSC モザイクキング・スタッキング





# HSC カタログ作成



# HSCデータ解析の手順(アルゴリズム) における課題

- (1) 大きく、かつ変動するDistortion
  - (2) 視野内の大きな大気差
  - (3) 広くて(CCD104枚)、円い視野、広いサーベイ面積
  - (4) 大きな蹴られ
- } ジオメトリ  
} 広視野



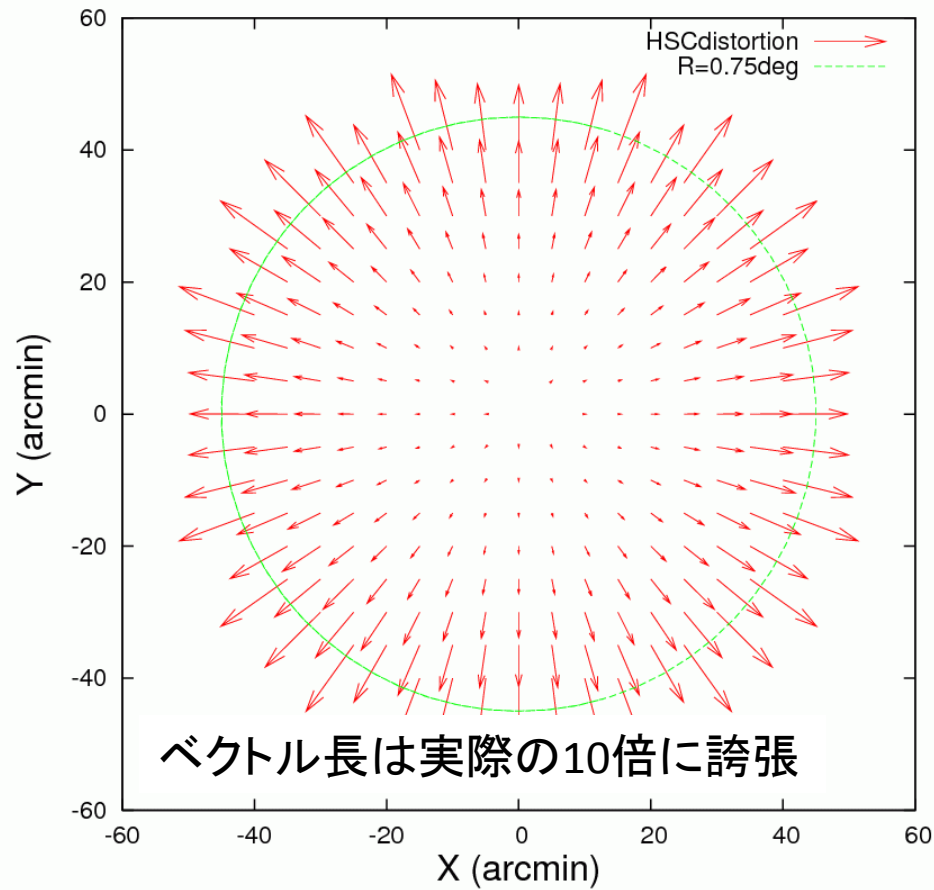
これらは相互に関連し、

1. 位置較正
  2. フラックス較正
  3. 天体の形状保存
  4. モザイクング・カタログ作成
- を困難にする

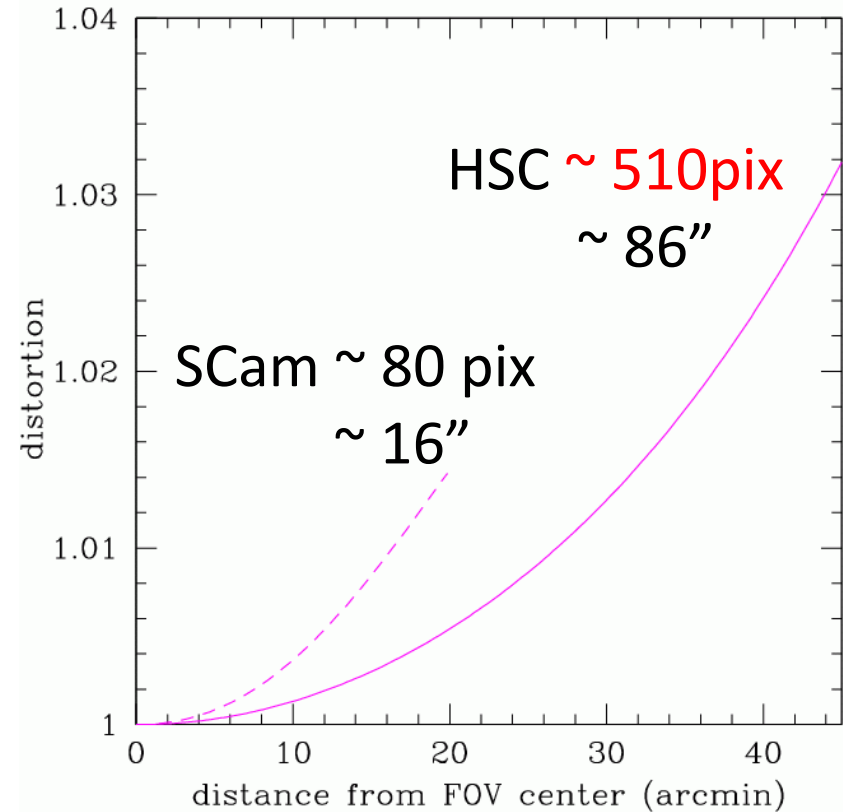
(1)(2) Distortion、大気差による困難

# 大きな Distortion

- 視野端ではSCamと比べて5倍歪む



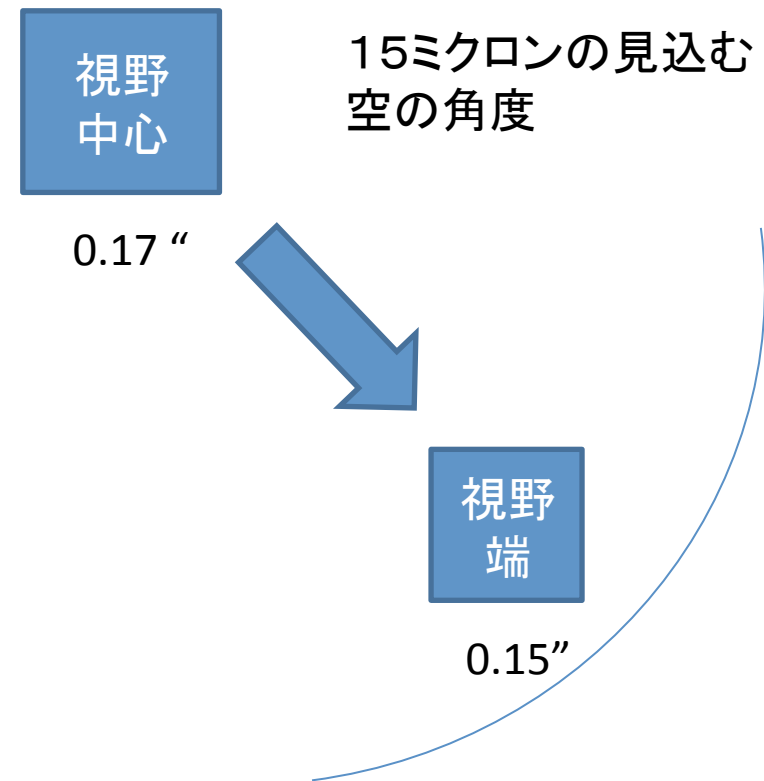
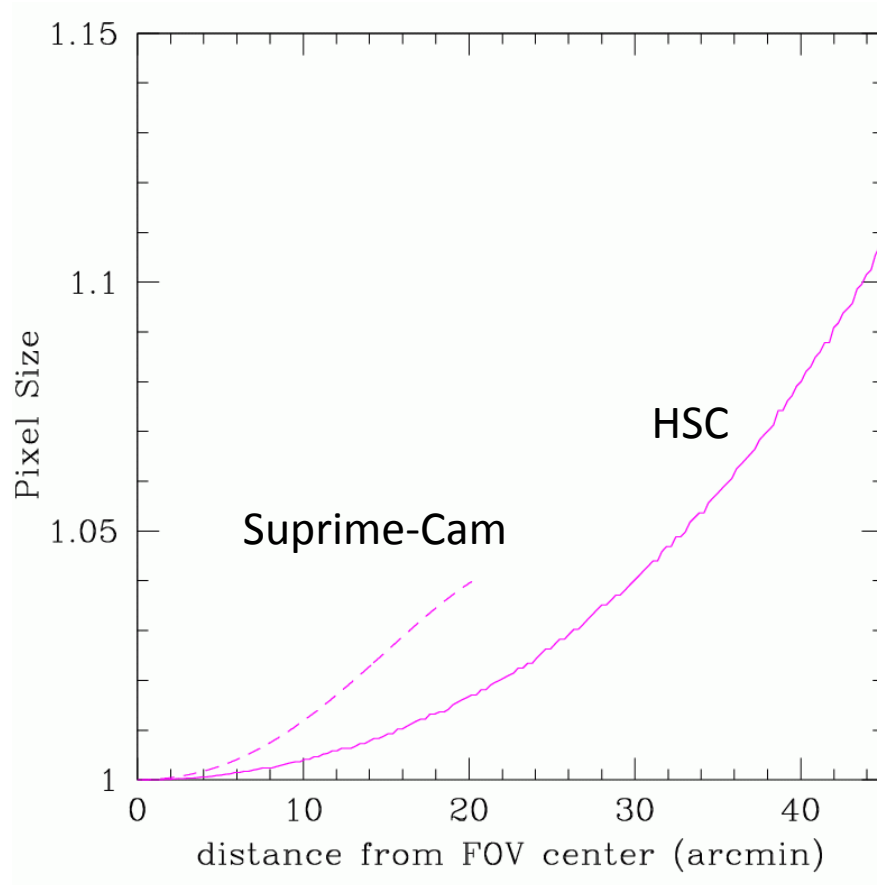
HSCのDistortion設計値



Thanks to 小宮山、諸隈、大倉

# 大きなDistortion

- 視野端でピクセルスケール(ゼロ点)差~11%
  - $\ll 0.05$  magレベルの測光で無視できない

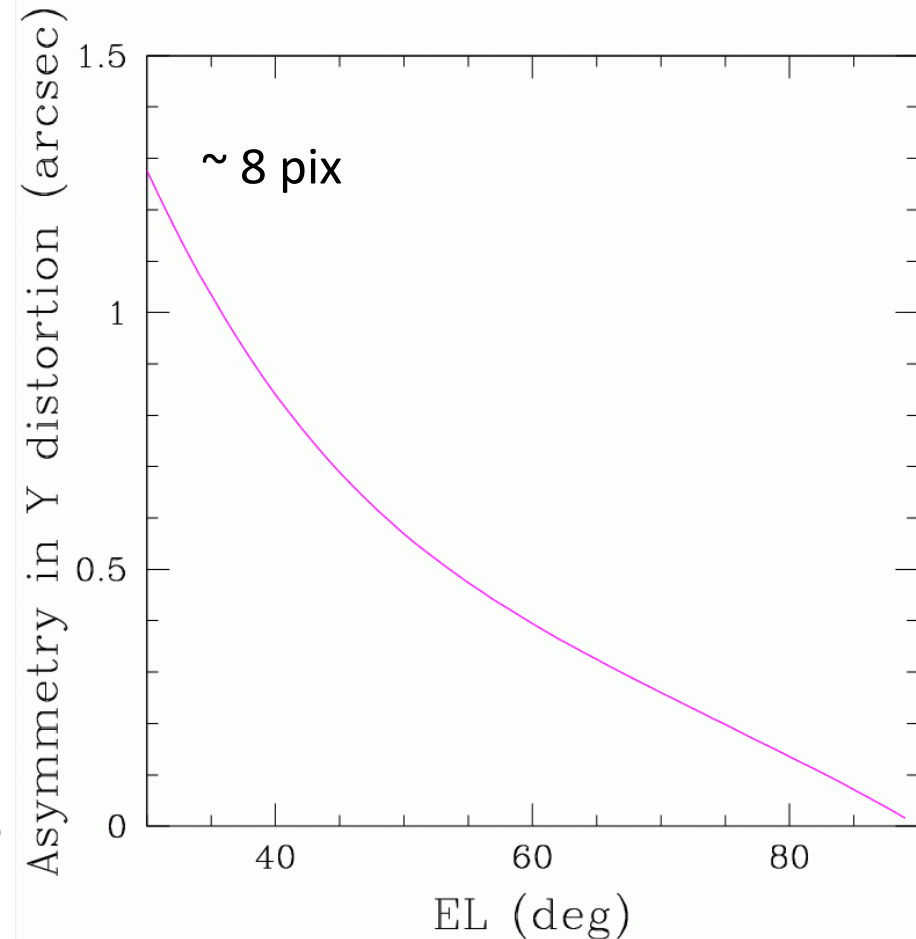


# 非対称で変動するDistortion

- ADC(大気分散補正光学系)の影響でDistortionパターンが、視野の重力方向に非対称性を持つ
- EL = 30 – 90 で  
~1.3 arcsec程度の変化

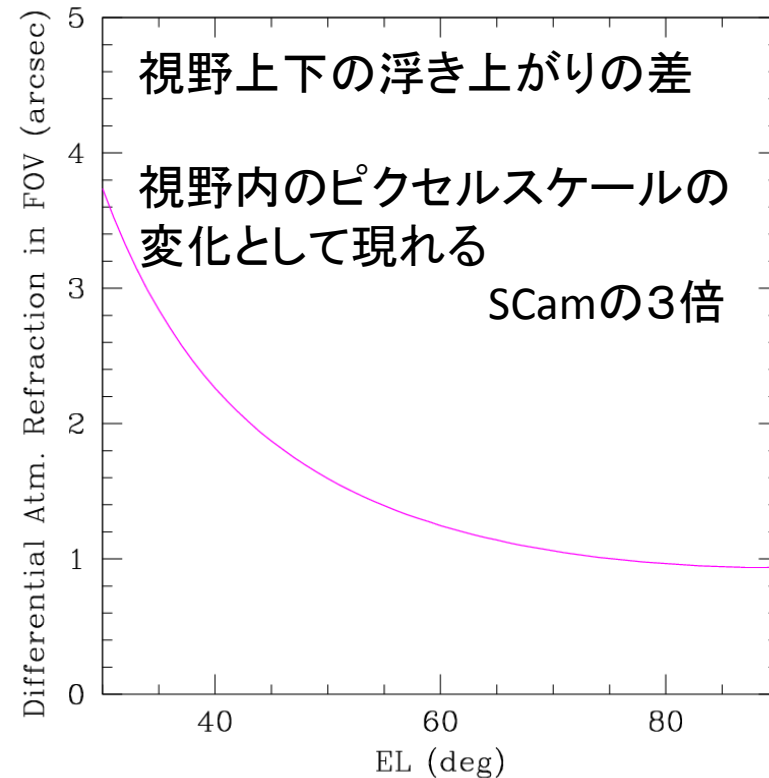
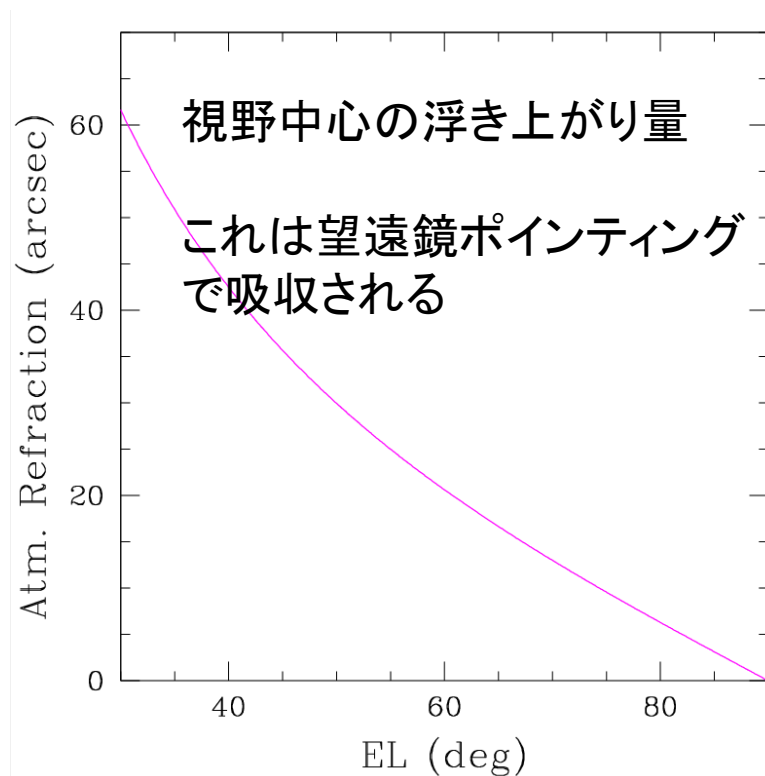


0.1pixオーダーでの  
形状測定、モザイクング  
にはDistortionを正しく決定  
しなければならない  
最長積分時間にも影響する



# 大きな微分大気差

- 視野上下の浮上量の差はEL = 30 – 80で2.5秒角  
–  $dr/dz = 0.''624 \text{ sec}^2 z (\text{deg}^{-1})$  (Tanaka 1993; ~V band)
- 視野が広いので実は水平方向に1"差がある



# 解決策の開発：ジオメトリ編

これまでの困難の影響をまとめると

- 位置・座標較正が非常に行いづらい
  - 参照星とのクロスID(マッチ)も困難
- 決めうちの補正では誤差も大きくなる
- ピクセルスケールがショットごと、視野内の場所ごとに変化するので
  - 天体形状の測定に系統誤差が入りやすい
  - 測光に場所に依存する系統誤差が入りやすい
  - 実は追尾にも影響する

SCamよりも強く、無視しづらい量の効果として現れる

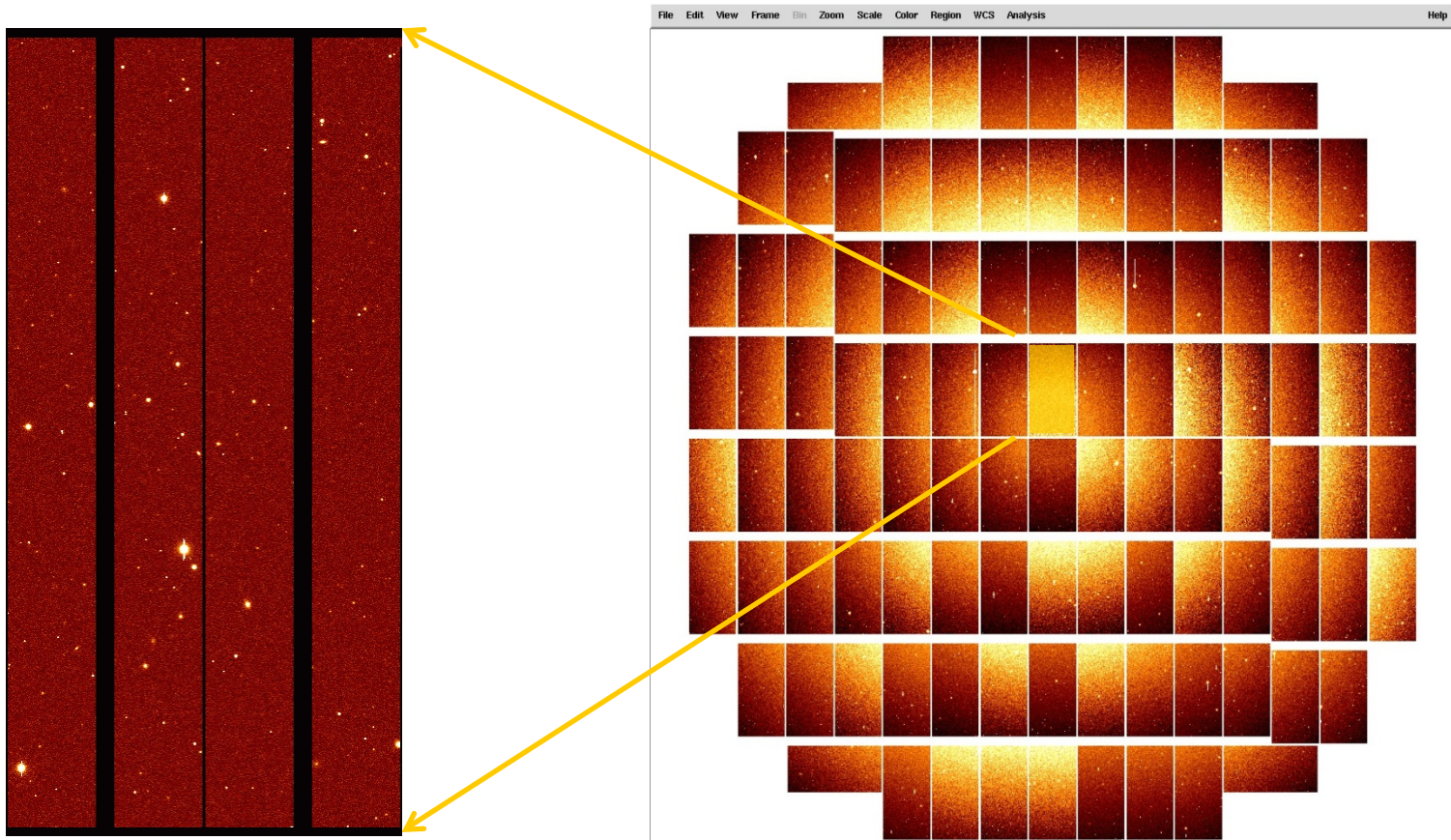


# 位置較正 (アストロメトリ) が重要

- 実際の参照天体を用いた、ショットごとに行う位置較正が非常に重要
- 後に述べるモザイクングのためにも重要
- 開発要素
  - 安定した参照天体と検出天体のカタログマッチング
    - ヘッダのWCSから  $> \sim 500$  pix のオフセット
    - CCD内の差分があり
  - 視野全体にわたる正確な座標変換式 (WCS; pixel – RA, Dec) の導出
    - 最終画像で各ショットごとに rms = 0.1pix オーダーの位置決定精度を目指す (仮: 絶対較正  $\sim 50$  mas, 相対  $< 20$  mas)

# HSCシミュレーションデータ

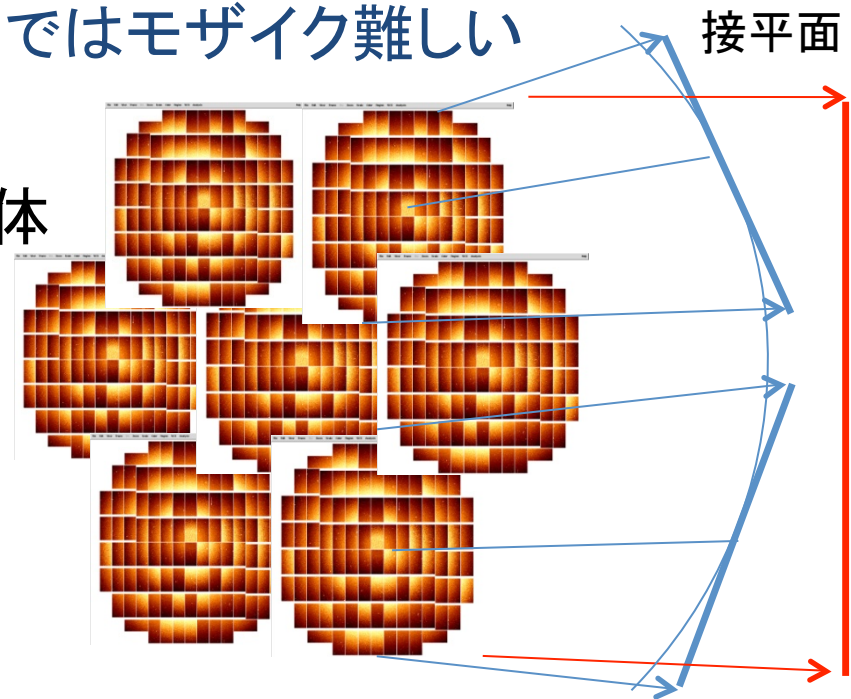
- 実際の観測データをシミュレートした104CCDを作成し、SCamの実データとともに、解析エンジンの開発に利用
- Distortion、大気による浮き上がり、けられを含む
- SDSSまたはUSNO-B1.0の参照星＋擬似銀河



(3)(4) 広く円い視野、けられによる  
困難

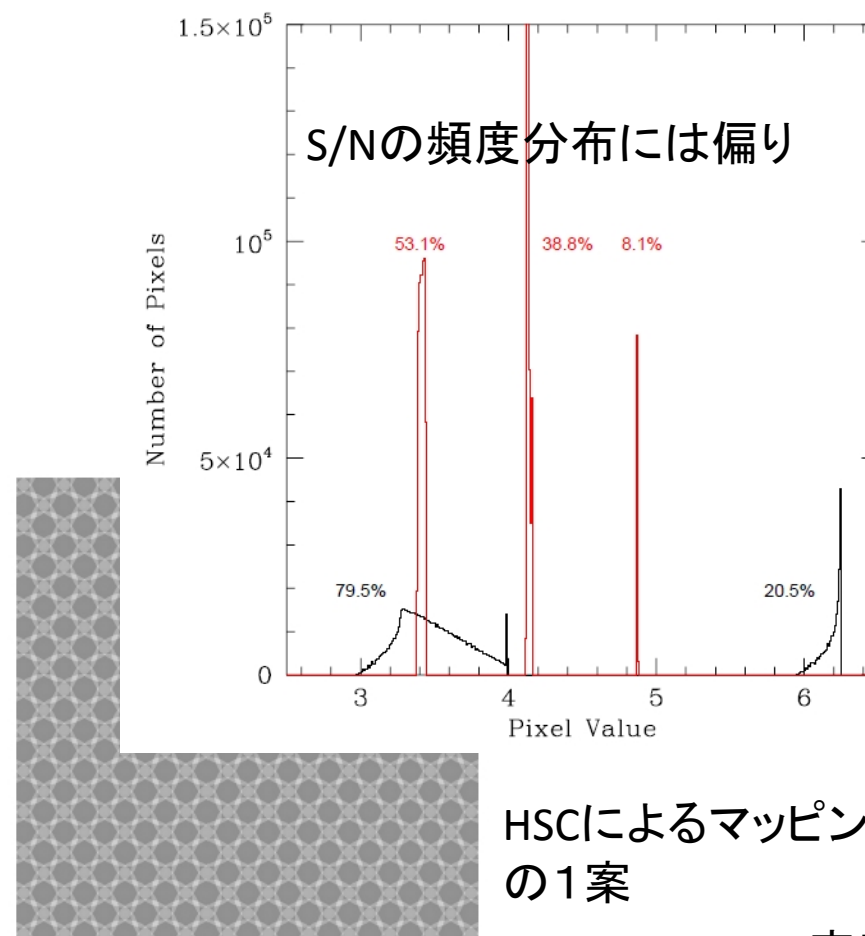
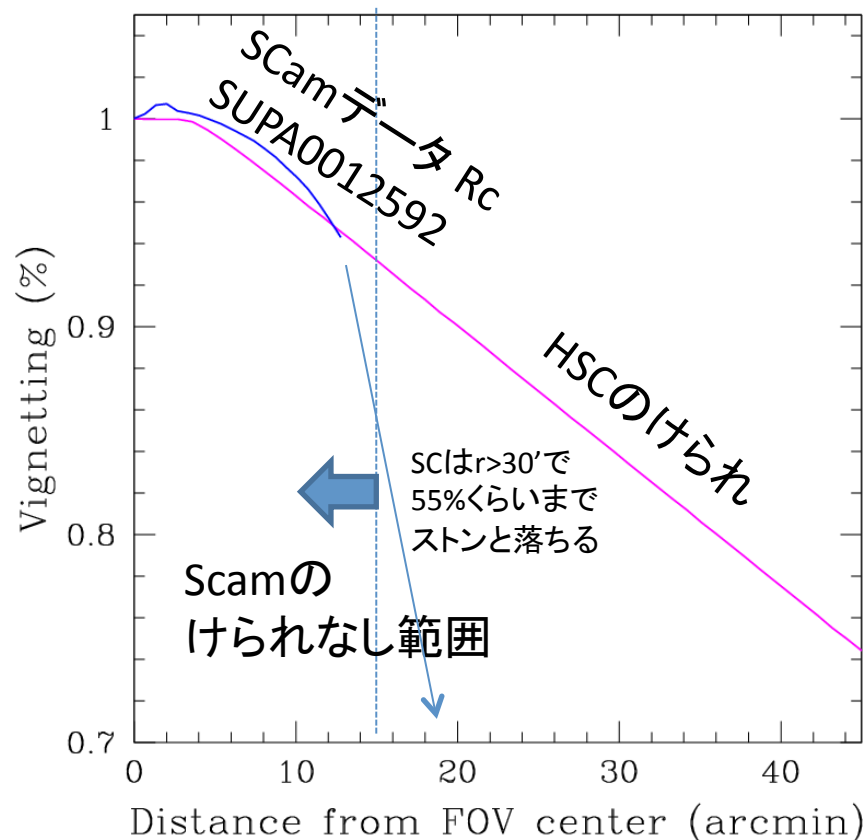
# 広く、かつ円い視野

- 数十～1000平方度＝HSCの10～>500視野でカバー
  - 各視野が重なり合いマップされることを想定
- どうモザイクするのか → 工夫が必要
  - 各視野は別々の天球の接平面
  - 仮にDistortion(接平面座標からのズレ)を除去できたとしても、平行移動と回転だけではモザイク難しい
  - しかも、実際は  
Distortionの除去(決定)自体  
が大きな課題



# 視野のけられ

- HSCでは視野中心から線形に落ちていく
- 広視野サーベイ → 複数ショットからのS/Nの異なるCCDの重ね合せ



HSCによるマッピングの1案

Thanks to 安田



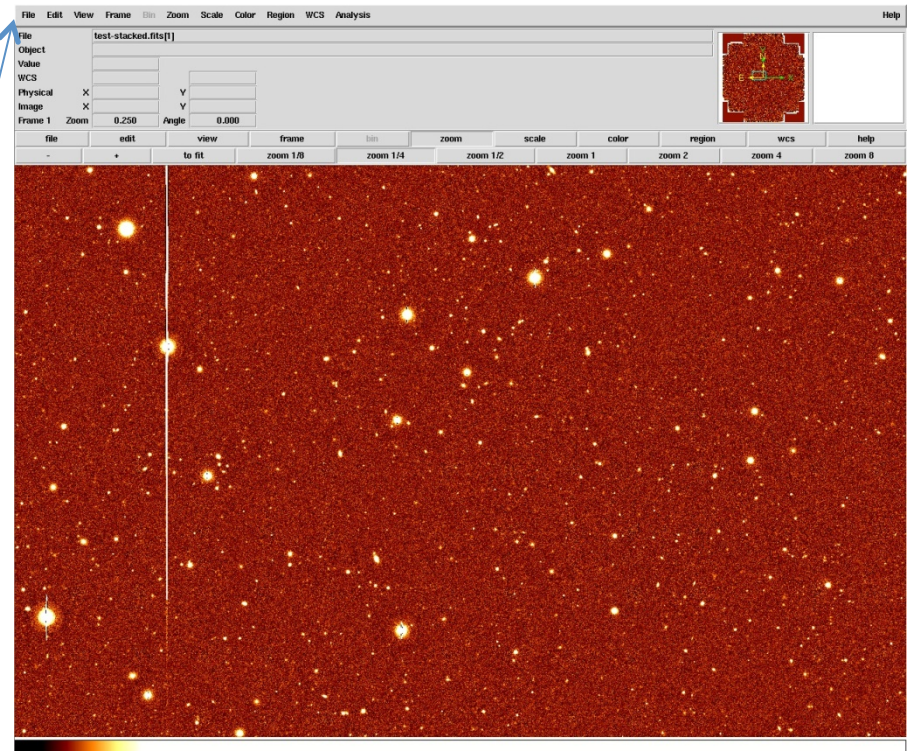
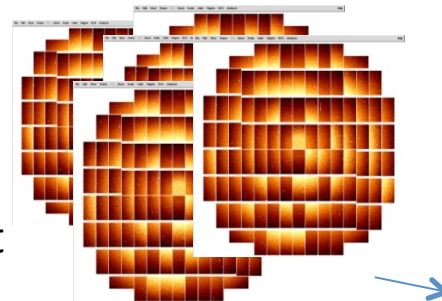
# モザイクキングの取り組み

- タイル=1視野として試験
- 共通天体の座標を使って、最小二乗法で各 CCDのTAN-SIP係数を決定  $\sim \text{rms} < 0.05 - 0.3 \text{ pix}$
- 課題

- モザイク精度の確保
- 視野をまたいだタイル
- 高速化
- ノイズ、マスク、PSF

データ諸元:

- EL=30
- SDSS stars + 偽銀河
- iバンド 300s x 5 shot



安田

# そのほかの解析手順の課題

- 場所ごとのPSFの測定と、それを用いた天体検出測定
- スタック時のPSFの均一化
- ピクセルごとのノイズ、マスク情報を用いた天体検出測定
- HSCに最適な天体測定アルゴリズムの確立
- HSCに最適な多色カタログ作成方法
- 正しいフラットをどのように作るか
  - Geometric effects, 散乱光・迷光etc. の影響の除去
- フラックス較正 (特に参照天体を含まないフィールド)
- 画像の座標変換とスタックをしないで、正確に形状測定する方法 → 宇宙論WLチーム

# 解析の効率化について



# データ解析の効率化における課題

大量 (SCamの10倍) のデータを迅速に処理しリリースするために

## 1. 解析に時間のかかる箇所の高速化

- ✓ ハードディスクへのファイル入出力を最小化
- ✓ ファイル入出力のトータル性能の最大化(と安定化)
- ✓ 分散、並列処理 — 例えば1次解析、モザイク・スタッキング

## 2. 解析履歴の管理

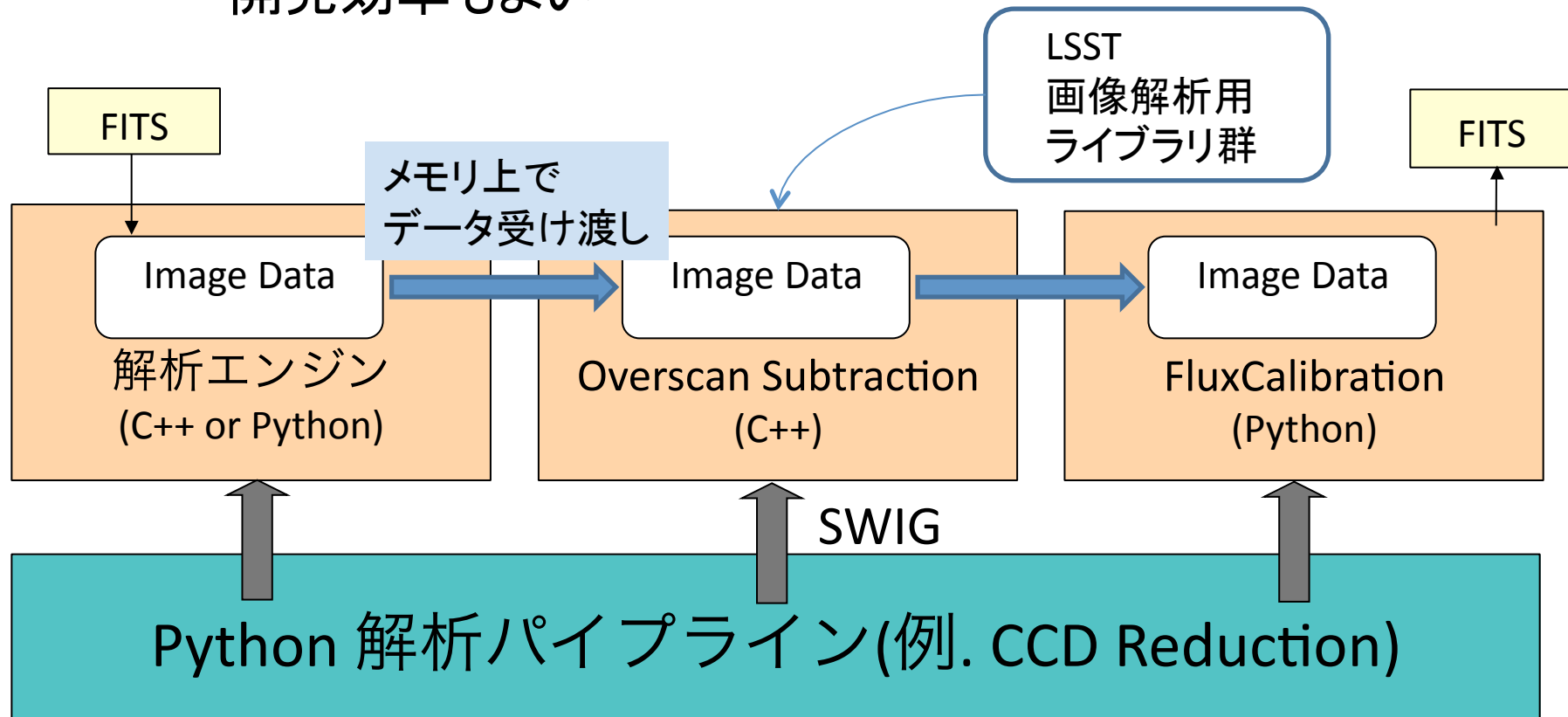
- ✓ 膨大な解析履歴と生成データを管理し、トラック・再現できるようにしたい

## 3. 観測時にデータクオリティを調べてタグ付けする機構

- ✓ オンサイト解析

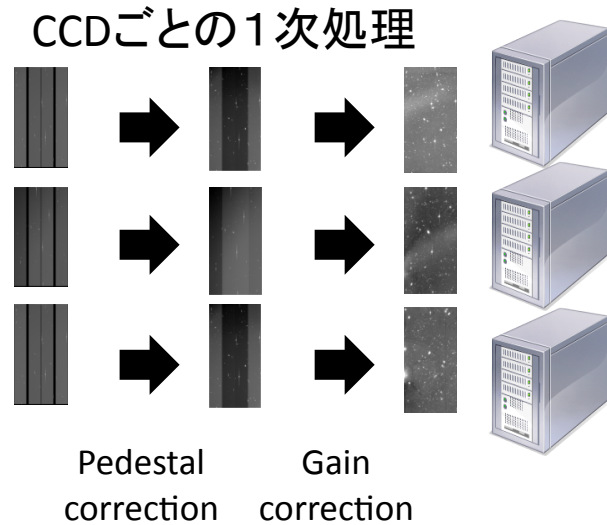
# 解析高速化： ディスクI/Oの最小化

- ハードディスクI/Oの最小化
  - 解析ステージ間はメモリ上でデータを受け渡す
- Python & C++(速度が重要なところ) + SWIG
  - 開発効率もよい



# 解析の高速化：並列処理の導入

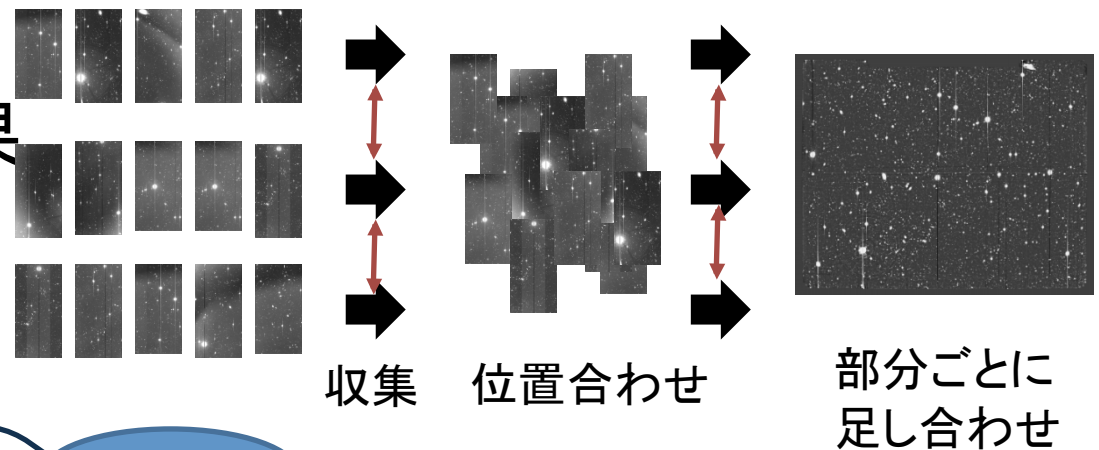
- データパラレル  
– 単純な並列化



複数台のコンピュータで  
たくさんプログラムを  
走らせるだけで  
並列化できる

- プログラムパラレル  
– 別の解析からの結果  
を集め、処理し、  
再分配（流れ作業）

CCDのモザイクキング

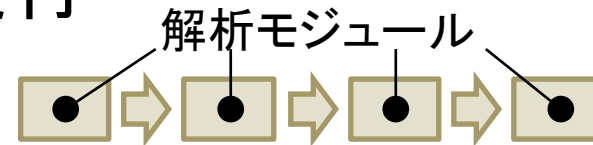


Thanks to 峯尾

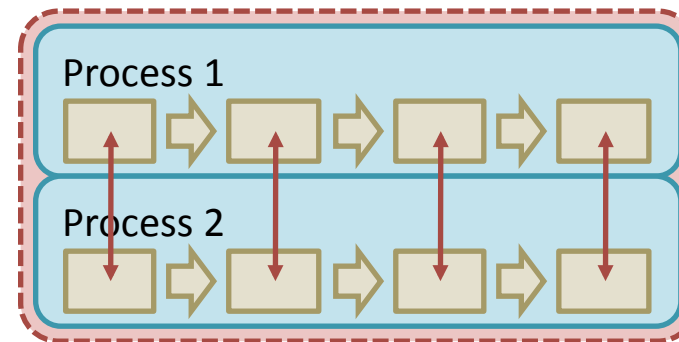
# 並列解析フレームワーク(pBASf)

[Belle Analysis Framework](#)

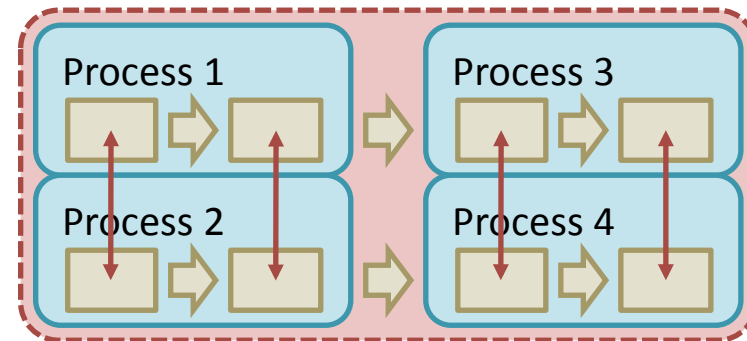
- 解析パイプラインの逐次実行
  - メモリを介したI/O



- MPIベース分散処理
  - データパラレル
  - プログラムパラレル
    - とともにプロセス間通信

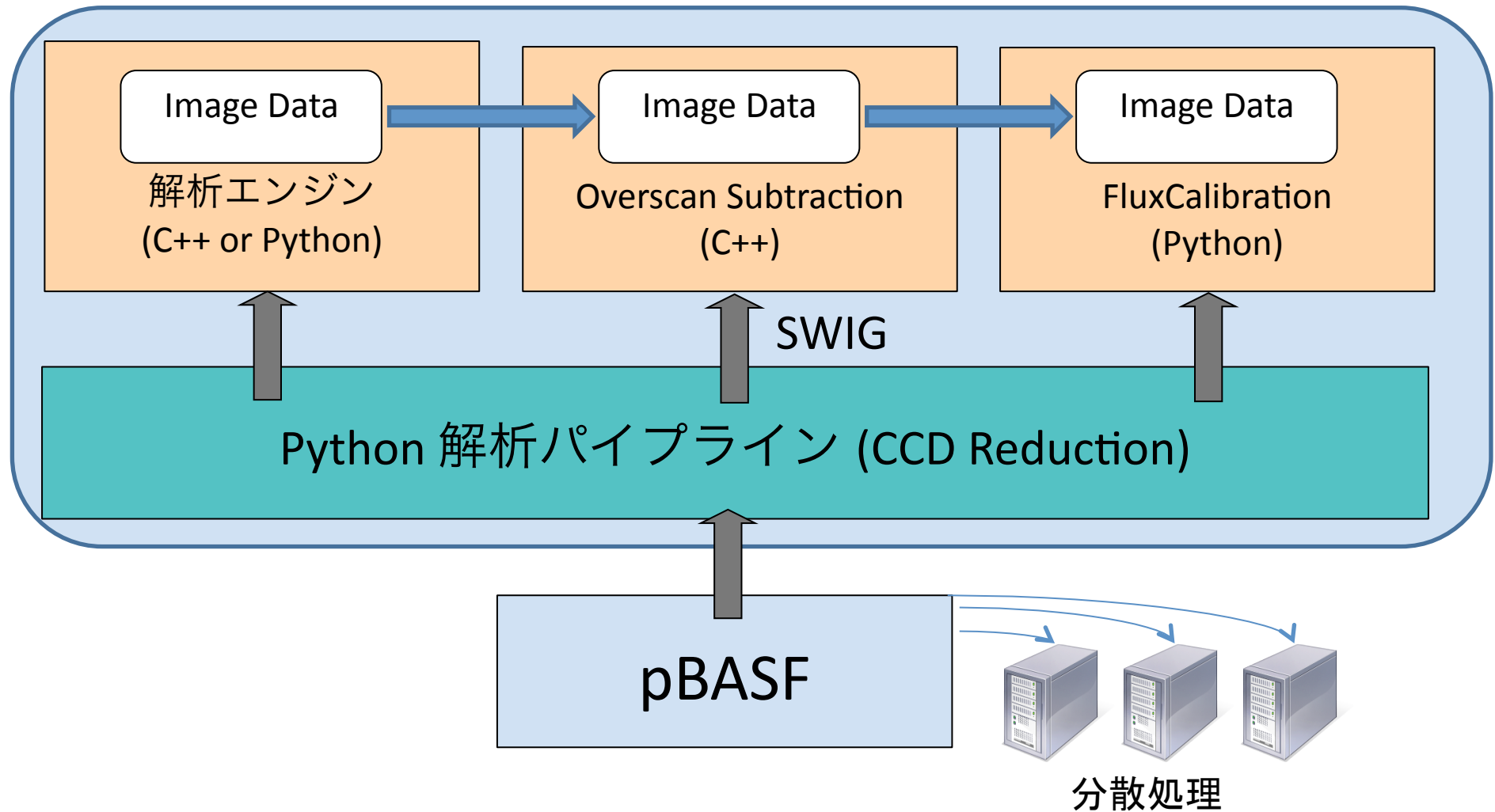


- Pythonインターフェース
  - Python、C++どちらで書かれた解析モジュールも呼べる



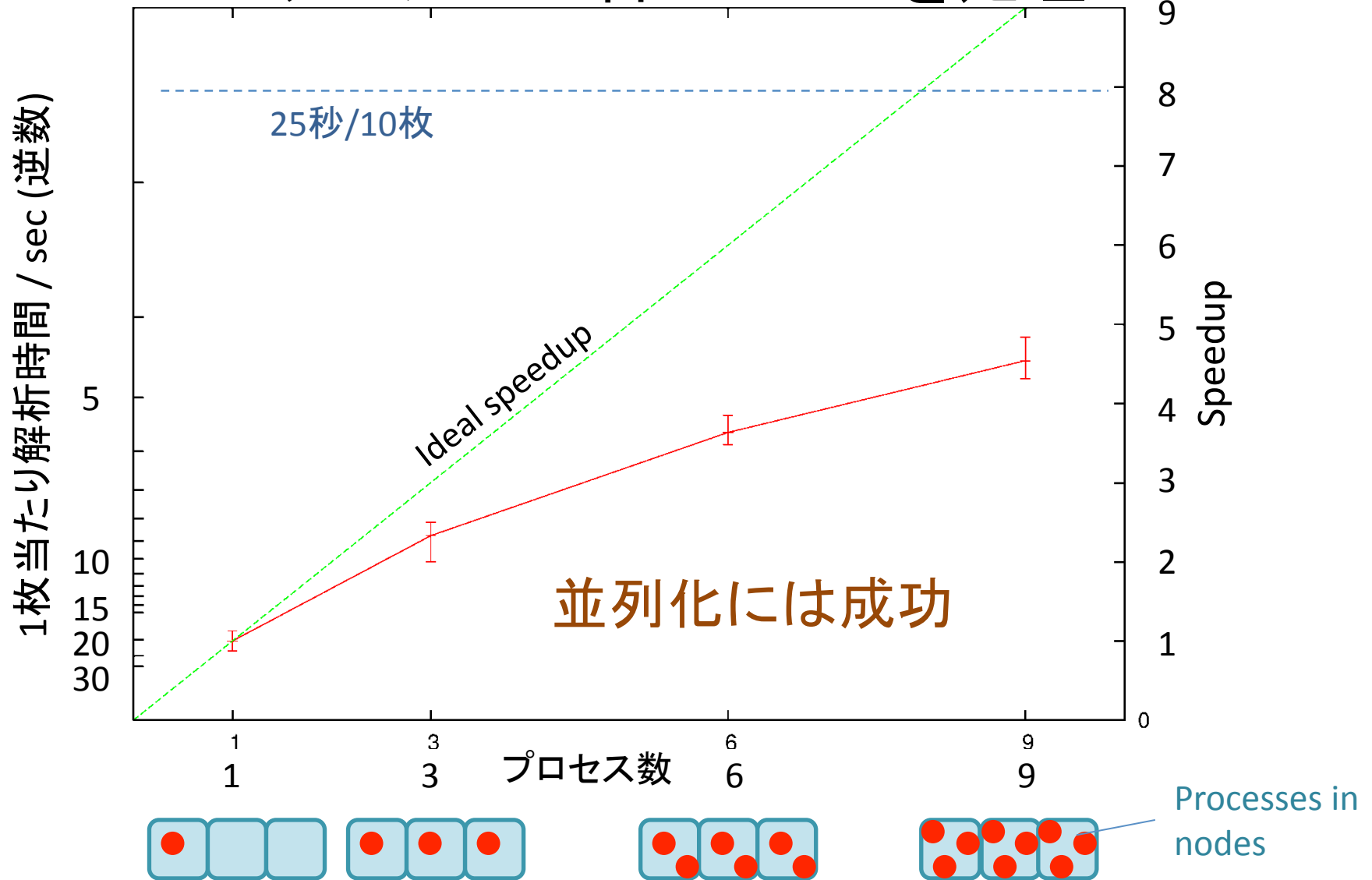
Thanks to 峯尾

# 解析パイプライン + pBASf



# Speedup

## 4コアマシンx3台で10CCDを処理



# 解析履歴の管理 & オンサイト解析

# 解析履歴の管理とオンサイト観測

## 1. 解析履歴の管理

長期の観測ランのサーベイデータの解析履歴とその生成データを管理し、トラック・再現できるようにしたい

» データアップデート時に的確な再解析・追加解析

## 2. オンサイト観測 = 観測時にデータクオリティを調べてタグ付けする機構

✓ 解析支援:

✓ 使えるデータかどうか、どのような素性のデータなのか、どのサーベイ領域のどのデータか(データセット)を記録。

→ フル解析で利用

必ず生データに付随させる

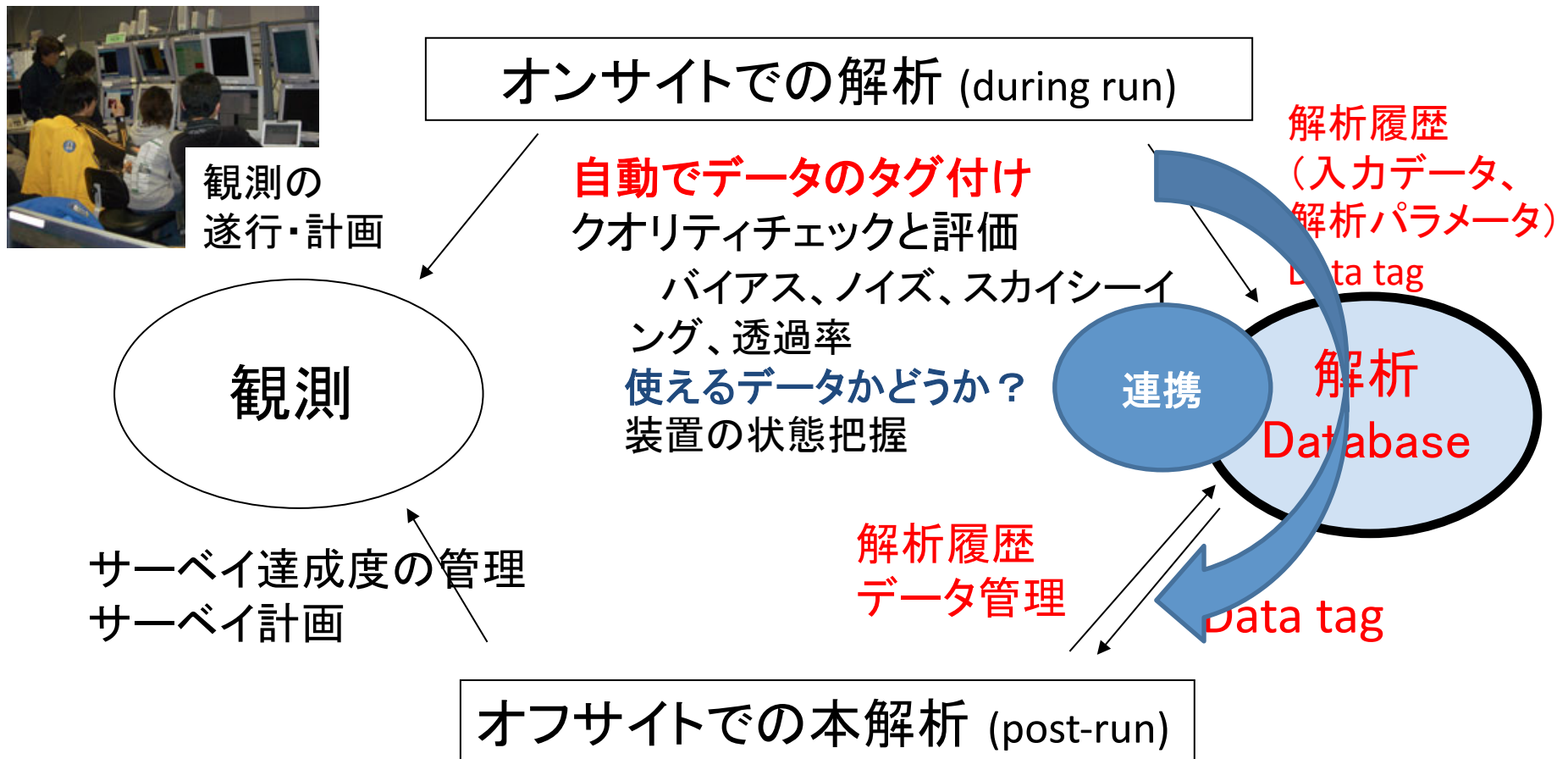
✓ 観測支援: サーベイ観測の柔軟な計画・遂行にフィードバック

解析履歴データベース共有



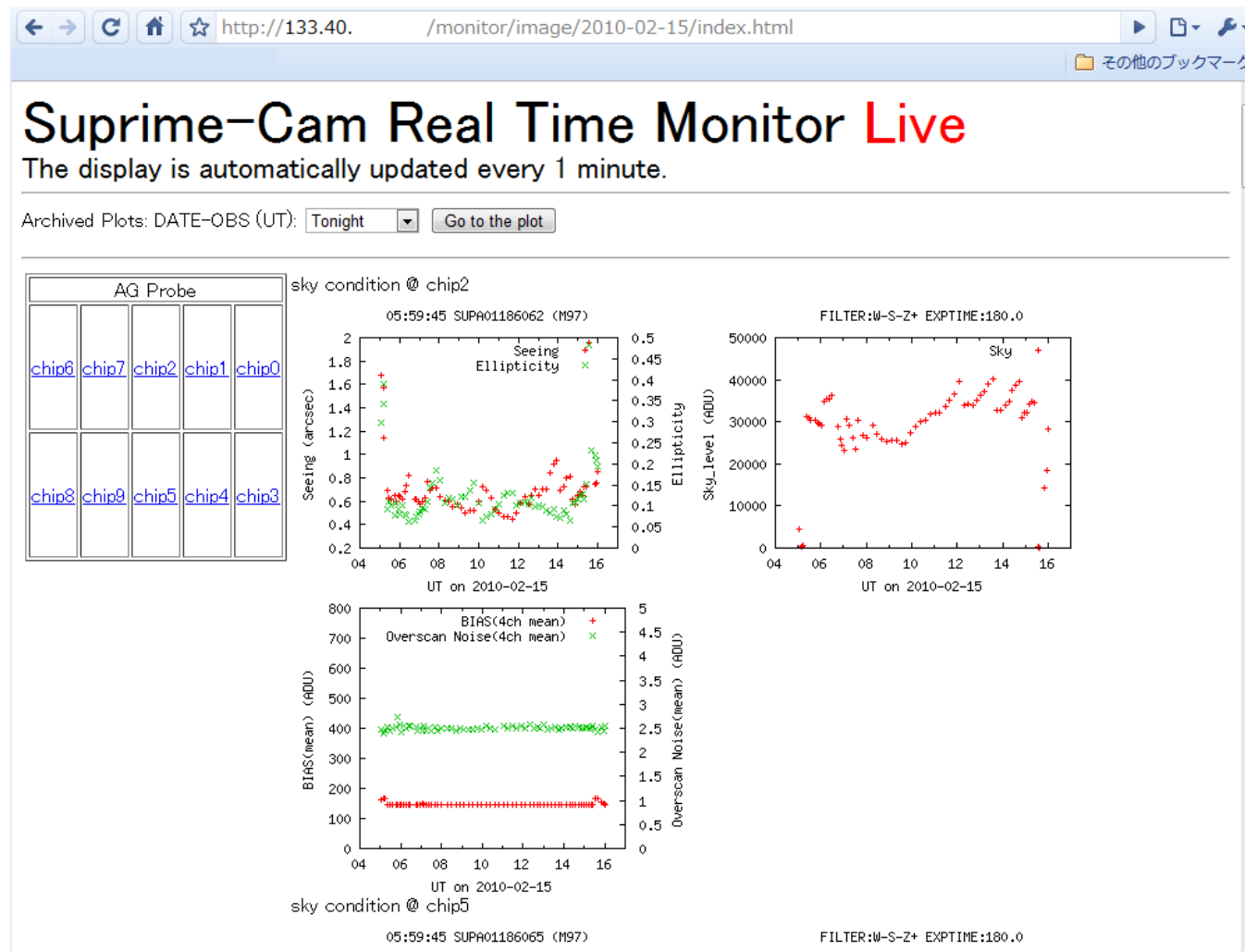
# データベースを用いた解析管理 オンサイト解析と本解析の連携

- ハワイ観測所にオンサイト解析システムのプロトタイプを構築し、**SCamの観測支援**と、解析の効率化のためのデータベース+ミドルウェアを試験。共同利用観測で運用試験している。



# 観測者用・時系列モニタ画面

- CFHT・Skyprobe、UKIRT・seeing monitorに類する機能の提供
- サーベイ管理ソフトウェアの開発も念頭に



# データベースの設計

(まだこれから、、)

# HSCサイエンスデータベースの役割

- 各研究者へHSCサーベイの最新の画像とカタログデータを提供
  - 主要なサイエンスに必要な基本情報を含むように設計
  - サーベイの進行に合わせて、一定期間ごとにアップデート&データリリースを想定
- フォローアップ観測準備のために必要な情報(位置、明るさなど)を提供
  - 既存のサーベイデータの情報ともリンク
  - 将来的な分光フォローアップなどに必須
- 複雑なQueryが可能(SQL文を直接入力)

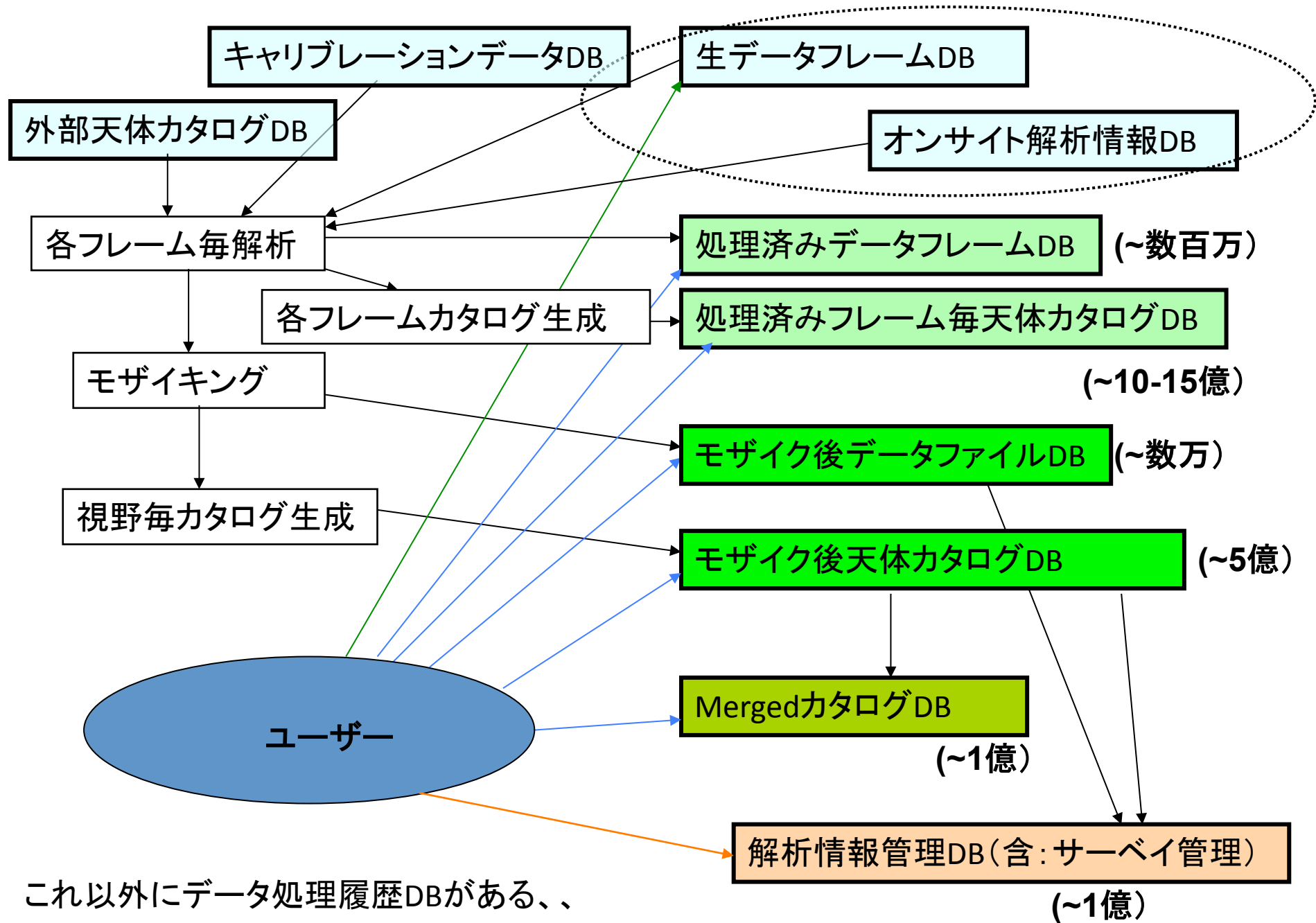
# データベース要件

RDBMSを使用予定 (PostgreSQL)

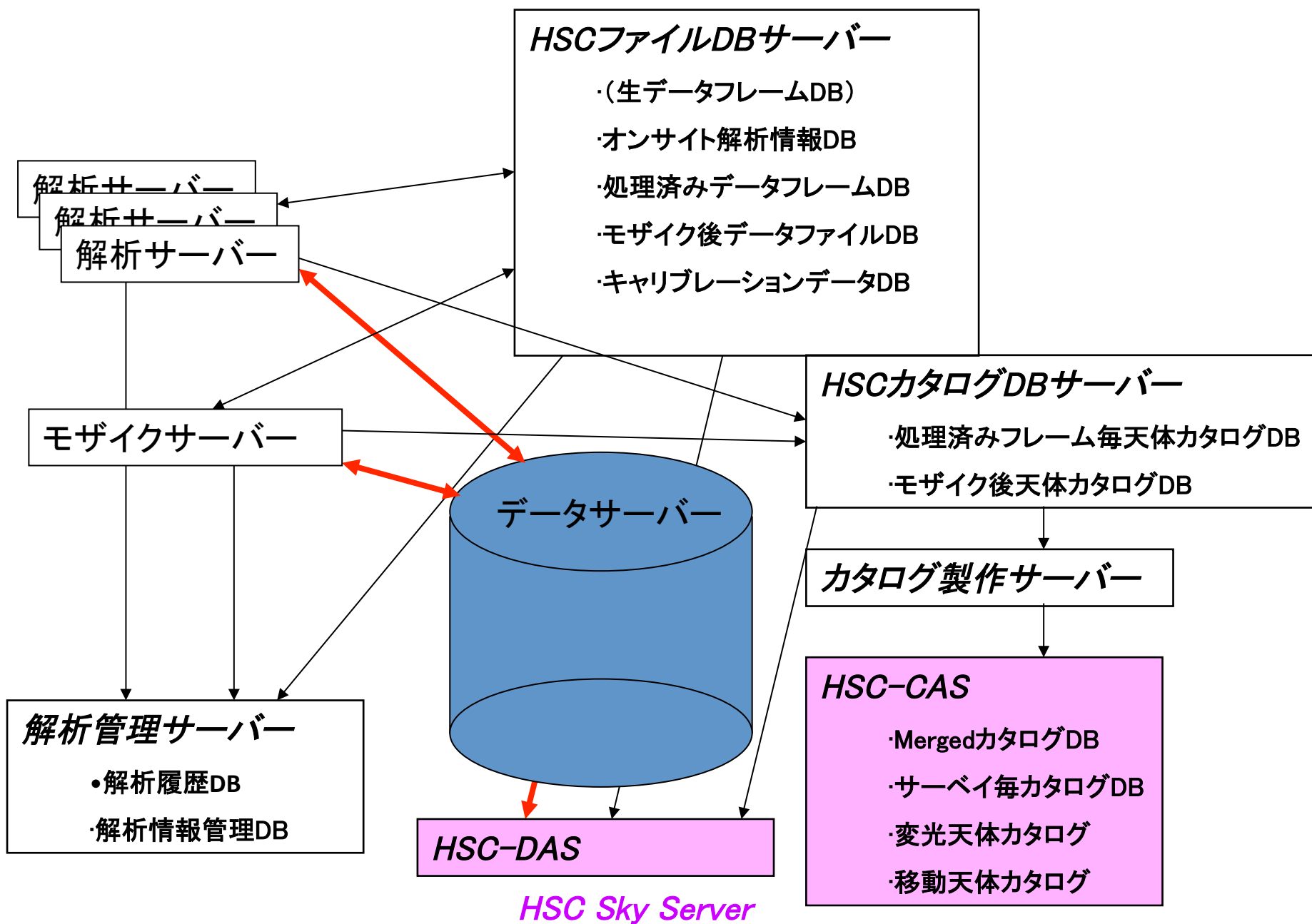
- 観測フレーム(FITSファイル)毎のレコード
- 最終画像での検出天体毎のレコード(最終天体カタログ)
- 各ショット毎(各フレーム毎)の検出天体カタログ
- ある程度の大きさの空間でサーベイの達成度を知る

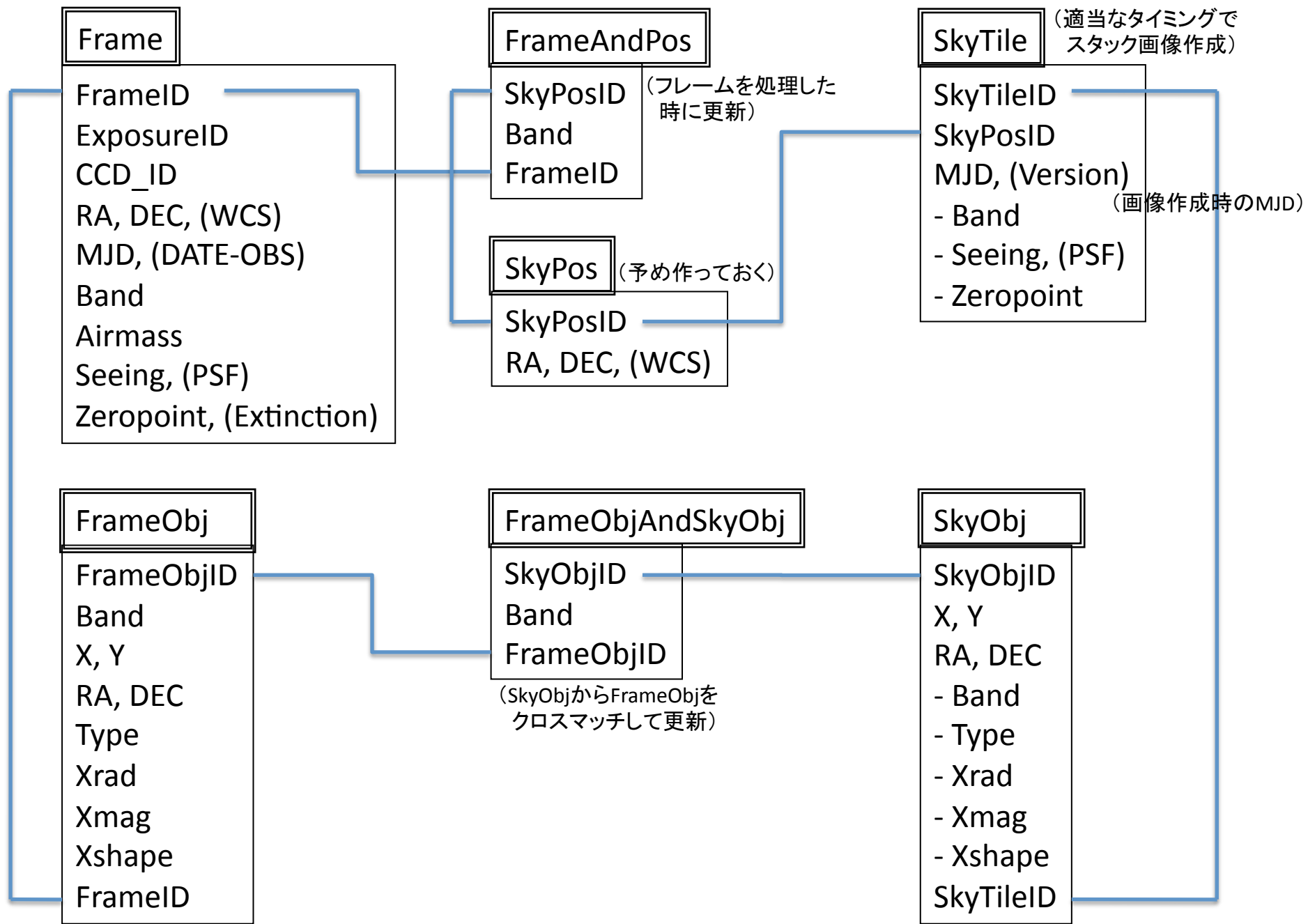
5年間で1000平方度を5フィルターで探査すると、、、

- ✓1露出104ファイル(合計で数百万レコード以上)
- ✓1000平方度で~26等までで~1億天体
- ✓1000平方度で26等まで、5フィルターで最低3回ずつ  
(合計で~15億レコード(天体)程度)
- ✓HealPixインデックス(1-2平方分程度)毎の達成度  
(1ショット6000インデックス × 数万ショット = 1億レコード)



# HSCデータ解析+データベースシステム構成模式図





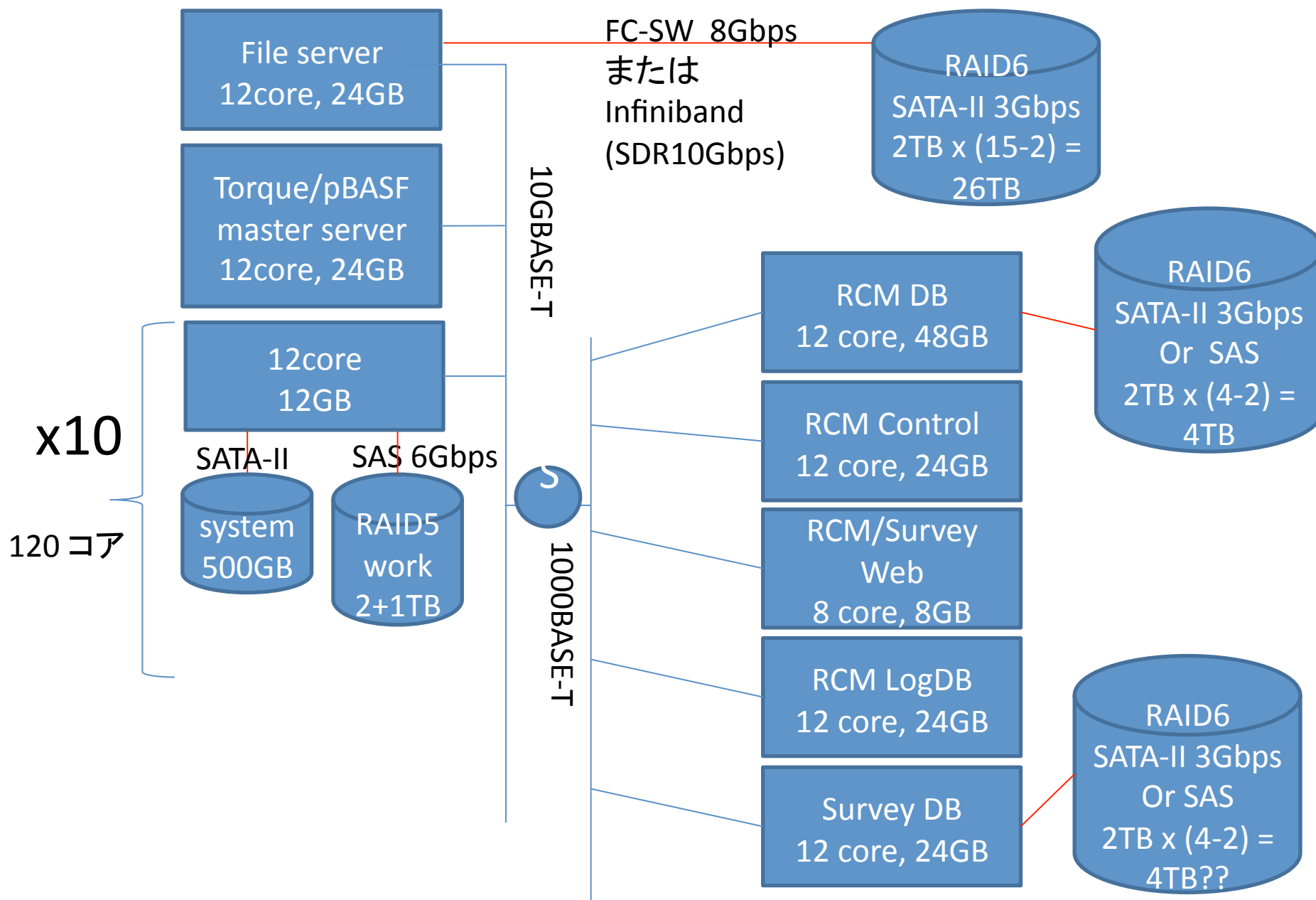


# 計算機システム

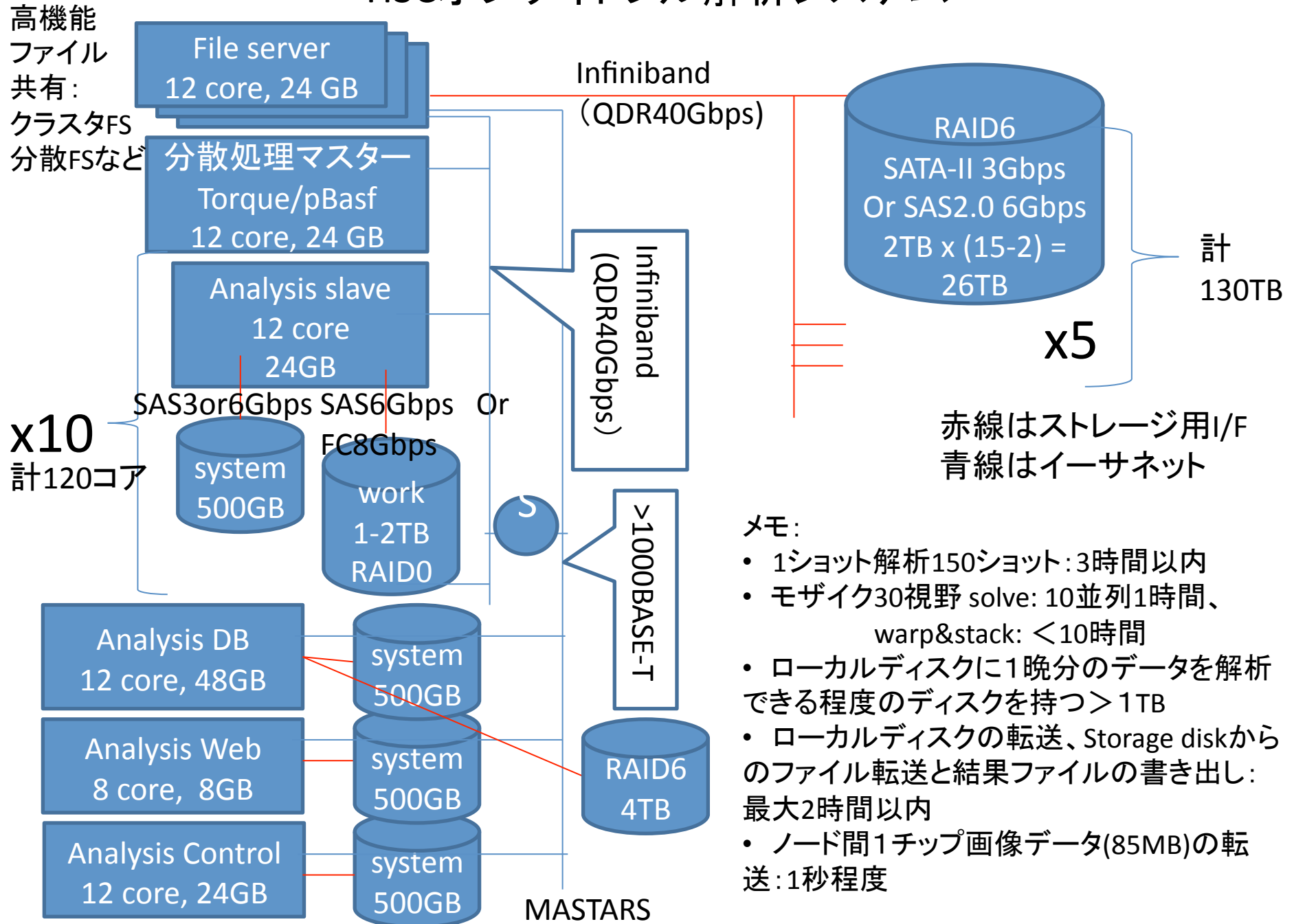
# 計算機システムの基本理念

- ◆なるべく安く、でも、安定した運用の確保
- ◆分散ファイルシステムを用いて高速データI/O確保  
(具体的には、NFSではなくLustreなどの導入を検討)
- ◆データ解析と観測運用とは一体だが、お互いの干渉はなるべく起こらないようにする。
- ◆解析システムは多人数のログインを想定せず
- ◆データベースには多人数同時アクセス

# オンサイト解析

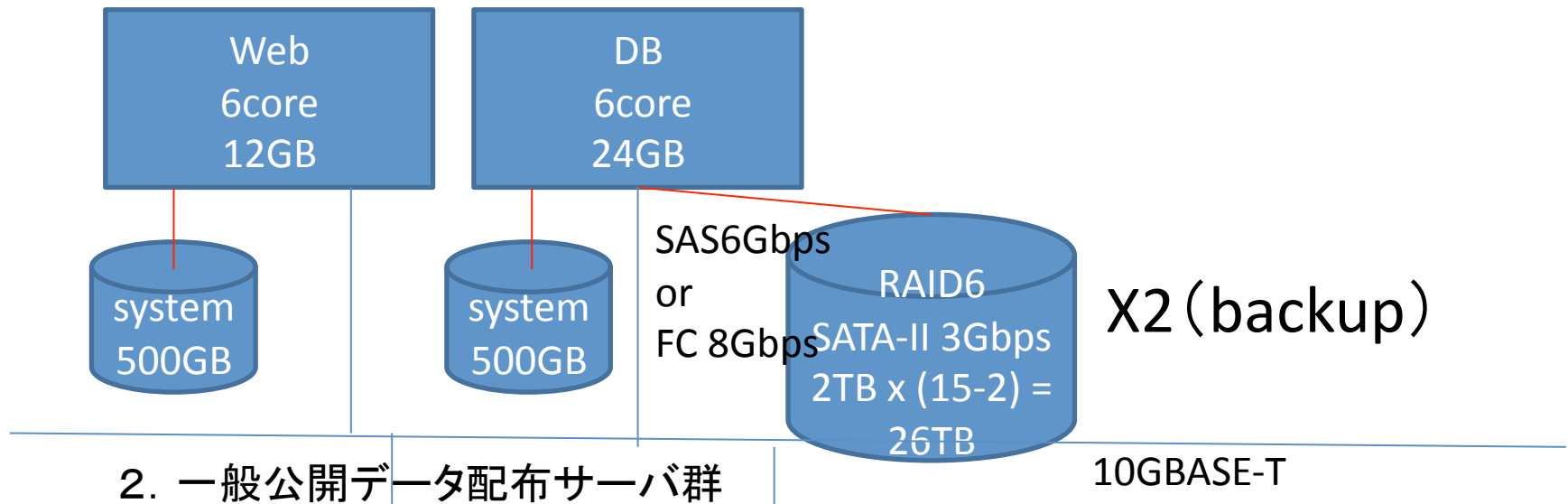


# HSCオフサイトフル解析システム

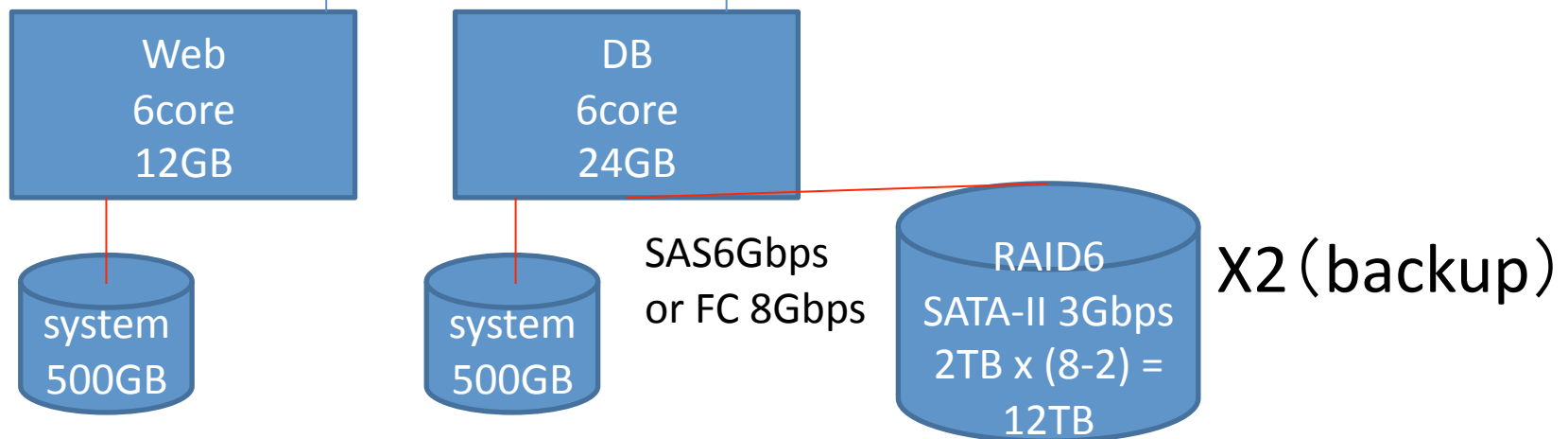


# データベースサーバ

## 1. サーベイメンバー用データ配布サーバ群



## 2. 一般公開データ配布サーバ群



# まとめ

- HSCはデータが量的には今まで(SCam)の10倍(1晩300GB)
- 画像データの扱いは今までと比べると随分複雑
- データフローを考慮して効率的な観測運用を目指す
- 画像処理とデータベースを統合した総合システム
- 分散ファイルシステム・並列化などによる効率化は必須
- サーベイデータを将来の基本データに、、
- 検索などの速度を確保するための検討はこれから本格化

**課題山積！！**