#### BH15.15/DDBJ

### DDBJ 登録システム周辺の開発

- Annotated Sequence登録に関する情報定義(真島、李、藤澤)
  - Structured\_Comment定義 <u>背景・導入と要約</u>
- DDBJ Validation Service開発
  - BioSample validationルール定義
  - o OWL/RDFを利用したBioSample validation
  - TogoAnnotatorを利用したproductのvalidation
  - Web Application開発

# BH15.15関連開発

- 微生物ゲノム配列解析パイプライン作成(鈴木治夫、板谷、横江、永田)
- 次世代シーケンサーのRDFデータスキーマの構築(理研:小林)
  - DDBJ DRAのXMLスキーマを参考にRDF化を試みる ←藤澤さんご議論お願いします
  - 理研作業フローの確認
  - o Excel(バイオロジスト) → RDF → DRA XML

## 1日目

## 目的

Structured COMMENTの入力情報についてユーザにDDBJ HPを通じてアウトプットする

- http://www.ddbj.nig.ac.jp/sub/comment-j.html
- http://www.ddbj.nig.ac.jp/sub/mss/annotation\_file-j.html
- #日本語版の差し替え: コンテンツは準備完了。
- #英語版の差し替え: コンテンツをなるべく早く準備。
- # VPN 接続に失敗した上に昨晩の停電 (3/16 20:45- 30min) で DDBJ の
- #機能停止中なので、戻ったらなるべく早く反映する。

# 作業

Structured\_CommentについてDDBJ定義情報としてスプレッドシートに整理する

- Structured COMMENT DDBJ
- <a href="http://www.ncbi.nlm.nih.gov/genbank/structuredcomment">http://www.ncbi.nlm.nih.gov/genbank/structuredcomment</a>

# 定義

- Structured Commentの表記は?"structured COMMENT"
- では、structured COMMENTとは?
  - INSDC (or RefSeq) の「flat file 上の COMMENT block において##[tagset\_id]-START## と##[tagset\_id]-END## の行に挟まれた行で "::" で区切られた tag\_name と tag\_value の組で構造化されたデータ記述。」
  - 記述される内容は、通常、その flat file に示される配列データ全体にかかる補足的説明であることが期待される。
  - o DDBJ annotation file では ST\_COMMENT feature
- "tagset\_id" とは
  - tag\_name, tag\_value を共有するデータセットの名称
  - DDBJ annotation file では ST\_COMMENT feature 配下の tagset id qualifier の値
- "tag\_name"
  - tagset\_id を共有する structured COMMENT において、共有されるデータ項目の名称

- DDBJ annotation file では ST\_COMMENT feature 配下の ユーザー定義 qualifier の名称
- "tag\_value"
  - tag\_name で示されたデータ項目に対応する個別の値
  - DDBJ annotation file では ST\_COMMENT feature 配下の ユーザー定義 qualifier に対応する value

## 作業方針

- まず、実の登録運用で INSDC 共有している以下のtagset\_idついて まとめる
  - GenomeAssembly-Data
  - Assembly-Data
  - International Barcode of Life (iBOL)Data
- その後に、NCBIで扱っているtagset\_idについても整理する
  - o MIxS
  - o HIV

### Todo

- <u>スプレッドシート</u>のdescription、description\_jaをすべて埋める
- exampleをすべて埋める
- 留意すべき点、制約、ルールなどをcommentカラムに記述する。
  - <u>Structured COMMENT DDBJ</u> 記載 の情報を反映させる
- INSDCリソースにおける関連するqualifier/attributeなどの情報を related resourcesカラムに記述する
- Valuesの全リストを現在のエントリーからパースする【Done】

# 微生物ゲノム配列解析パイプライン情報共有(鈴木、藤澤)

- ゲノム微生物学会「シアノバクテリアに特化したCyanoBaseゲノムアノ テーションパイプラインの構築」共有
  - メタデータ追加部分について、Prokka4ddbjコードを共有
- 他に、慶応ではゲノムアノテーション出力については inference/db\_xrefでアノテーションを足したい
- Product nameの大文字/小文字変換問題について確認
- その他、共通リソースがあれば共有していく

# TogoAnnotator辞書開発(山本、李、真島、藤澤)

# 辞書開発方針の確認

- before/afterのアノテーションセットとしては、以下の2つを準備する
  - o INSDC-UniProt
  - UniProt-RefSeq
- UniProtのReviewed/UnReviewed および RefSeqのstatusで辞書 データセットにフィルターをかけられるようにする
- RefSegのstatus情報について調査する(Todo: 李)
- 辞書データからtaxon単位でのフィルターも設計上可能
  - アノテーションデータを分けたいtaxonをリストアップできるとよい
  - まずはCyanobacteria(taxid:1117)で試行

### 2日目

# tagset\_id、tag\_valueの使用状況の確認

Structured COMMENT情報をDDBJ release 103 およびGenBank 210 wgs master エントリーからパースし

tagset\_idの頻度情報およびDDBJ submissionにおいて対応している Genome-Assembly-Data, Assembly-Data, International Barcode of Life (iBOL)Dataに限定してtag\_valueの頻度情報を算出し、スプレッドシートに追加した。

# RefSeq COMMENT分析。

ただし、NCBI 側でも document 整備が不十分な様子なので、憶測が入っているかもしれない。

# Structured COMMENT - RefSeq

↑ RefSeq status情報の調査(李) をマージ

### 3日目

#### TODO

- Sequencing Technology, Assembly Methodの入力具体例をリストアップ
- DDBJ 用にStructured COMMENTのドキュメント策定する(真島、李)
  - 技術文書的な structured COMMENT とその要素の定義
    - Annotated Sequence/Structured COMMENTのバリデーションルール策定
  - 配列データ登録用の annotation file への記載方法: ST\_COMMENT feature

- フラットファイルの読み方の説明: structured COMMENT
- ラップアップ用資料策定

#### 技術文書的な structured COMMENT とその要素の定義 (新規)

#どこかで正式に公開すべきかもしれないが、場所など未定

- 表記: structured COMMENT
- "structured COMMENT"
  - INSDC (and RefSeq) の「flat file 上の COMMENT block において ##[tagset\_id]-START## と##[tagset\_id]-END## の行に挟まれた行で "::" で区 切られた tag\_name と tag\_value の組で構造化されたデータ記述。」
  - ユーザー(データ登録者を含む)コミュニティが定義したデータセットを flat file を介して公開・共有することを可能にする。
  - ただし、実際には、INSDC 側の都合で記述させていることも多い。
  - 記述される内容は、通常、その flat file に示される配列データ (エントリ) 全体に かかる補足的説明であることが期待される。
  - o DDBJ annotation file では ST COMMENT feature の配下に記載される。

## "tagset\_id"

- tag name, tag value を共有するデータセットの名称。
- o DDBJ annotation file では ST\_COMMENT feature 配下のtagset\_id qualifier の値。

## • "tag\_name"

- tagset\_id を共有する structured COMMENT において、共有されるデータ項目 の名称。
- Flat file 上では COMMENT block の##[tagset\_id]-START## と ##[tagset\_id]-END## の行に挟まれた行内で "::" で区切られた左側に表示される。
- Tag\_nameを複数回、記載することはしない。tag\_value 相当の値が複数ある場合は ";" 区切りで列挙する。
- DDBJ annotation file では ST\_COMMENT feature 配下のユーザー定義 qualifier の名称。

#### "tag\_value"

- Tagset\_id と tag\_name で示されたデータ項目に対応する個別の値
- 値が複数の場合は ";" で区切って列挙する。tag\_nameを複数回、記載することはしない。
- Flat file 上では COMMENT block の##[tagset\_id]-START## と ##[tagset\_id]-END## の行に挟まれた行内で "::" で区切られた右側に表示される。
- DDBJ annotation file では ST\_COMMENT feature 配下のユーザー定義 qualifier に対応する value

#### ● COMMENT としての見え方

- COMMENT は 80 字で改行する。
- "::" で区切る位置は固定長ではなく、データセット内で最も長い tag\_name に合わせる可変長。
- Tag\_name が長い場合、[tag\_name] + ":: " + [tag\_value] で 80 字に達した際に 改行するが、"::" の後までスペースインデントする。

#### 配列データ登録時のannotation file Feature/Qualifier としての記載ルール

- 1つの ST\_COMMENT featureに対して、1つ tagset\_id qualifier の記載が必須
- 1つの entry に ST COMMENT feature を複数回、記載して良い。

- 1つの entry に同一の tagset\_id を持つ ST\_COMMENT feature を複数回、記載することはできない。
- 1つのST\_COMMENT feature 配下に同一の tag\_name (ユーザー定義 qualifier) を複数回、記載することはできない。
- tag\_value で示すべき値が複数ある場合は ";" で区切り、列挙する。
- いくつかの tagset\_id の値と tag\_name (qualifier 名) は辞書化してチェック (jParser) する予定。3月末実装。

下記既存サイトの書き換え: 2016.03.18 済

配列データ登録用の annotation file への記載方法: ST\_COMMENT feature

http://www.ddbj.nig.ac.jp/sub/mss/annotation\_file-j.html#5-4 http://www.ddbj.nig.ac.jp/sub/mss/annotation\_file-e.html#5-4

Structured COMMENT 説明ドキュメント http://www.ddbj.nig.ac.jp/sub/comment-j.html

下記既存サイトの書き換え: 2016.03.22 済 Structured COMMENT 説明ドキュメント http://www.ddbj.nig.ac.jp/sub/comment-e.html