

# スパコンを活用した がんゲノム配列データ 解析手法の開発

愛知県がんセンター  
システム解析学分野



# はじめに

- がん細胞は、正常な細胞中の**DNA (ゲノム) 配列**に変異が起きることにより、細胞が異常に増殖する能力などを獲得したものです
- 同種類のがんであっても、ゲノム上で変異の起きる場所や種類の頻度にバリエーションがあります
- 現在**ゲノム上の変異**に応じて治療法の選択を考える、**がんゲノム医療**が、本格化しつつあります
- ここではがんゲノムの変異を見つけるための**元となる観測データ**はどのようなものか、**どのような原理で変異を検出している**のか、また実際にはどのような難しさがあるのかなどを、システム解析学分野での研究を交えて概説します

# システム解析学分野について（1）

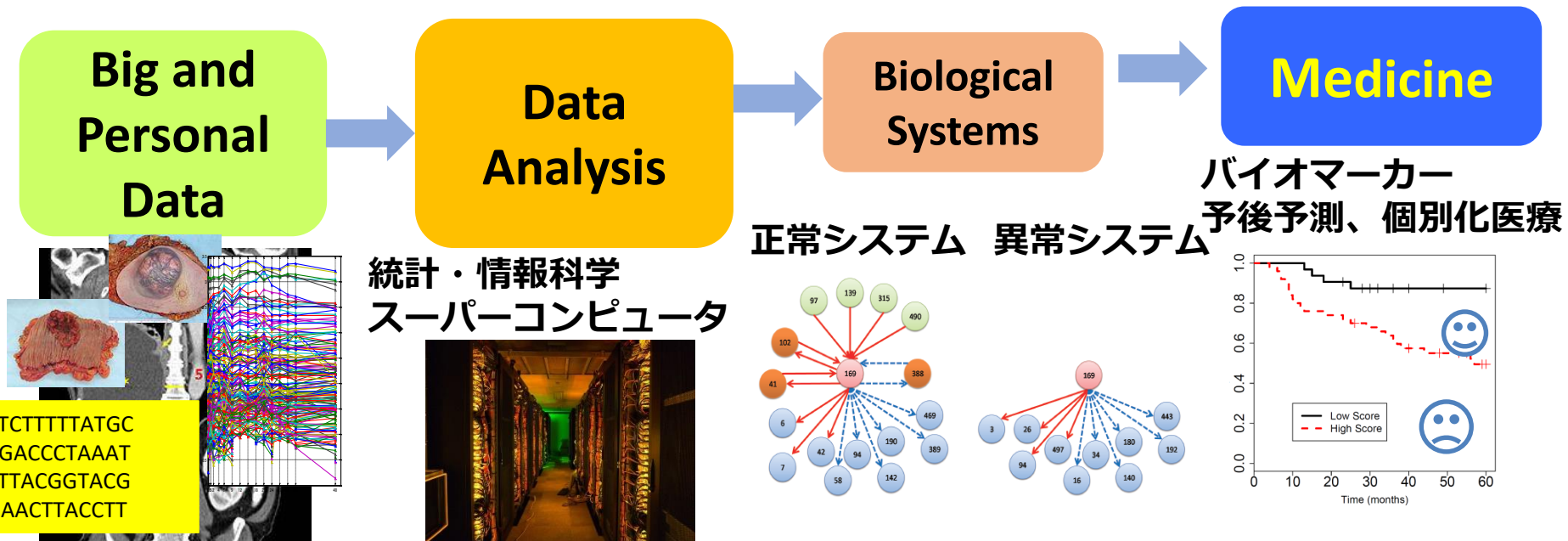
- 2019年2月に研究所にできた新しい分野です
- がん細胞から得られた**ゲノムデータ**などの様々な**生体ビッグデータ**を**解析する方法**の研究を行っています
- そして実際に、患者さんのデータから、**がん細胞の複雑なシステム**に関わる情報を抽出し、それを基に**一人ひとりに合わせた医療**へつなげることを目指しています



# システム解析学分野について (2)

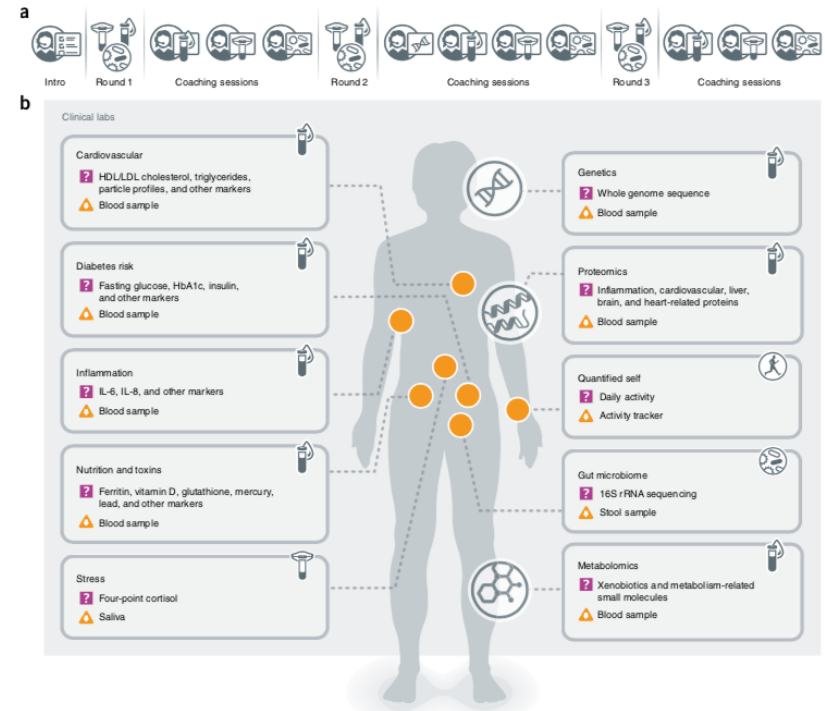
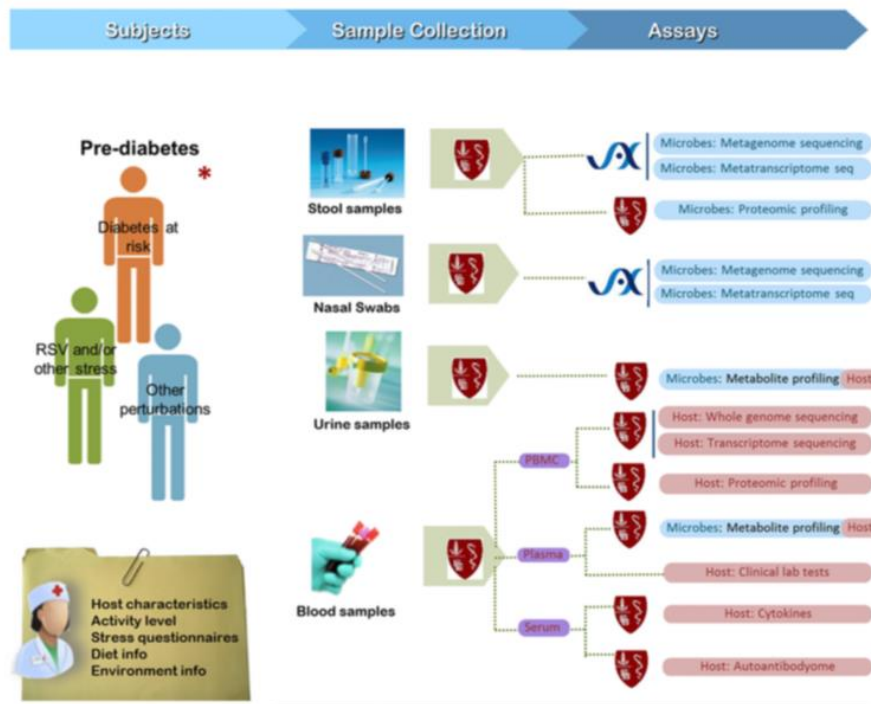
## • 研究テーマ

- メディカルバイオインフォマティクス
- 大規模生体データからの有用情報抽出のための情報科学・統計科学的データ解析手法開発
- 人工知能を活用した個別化医療のための情報解析基盤開発



# 近年様々な生体データの取得が可能になってきた DNA, RNA, タンパク、メタゲノム、、、

Integrative Human Microbiome Project      Pioneer 100 wellness project

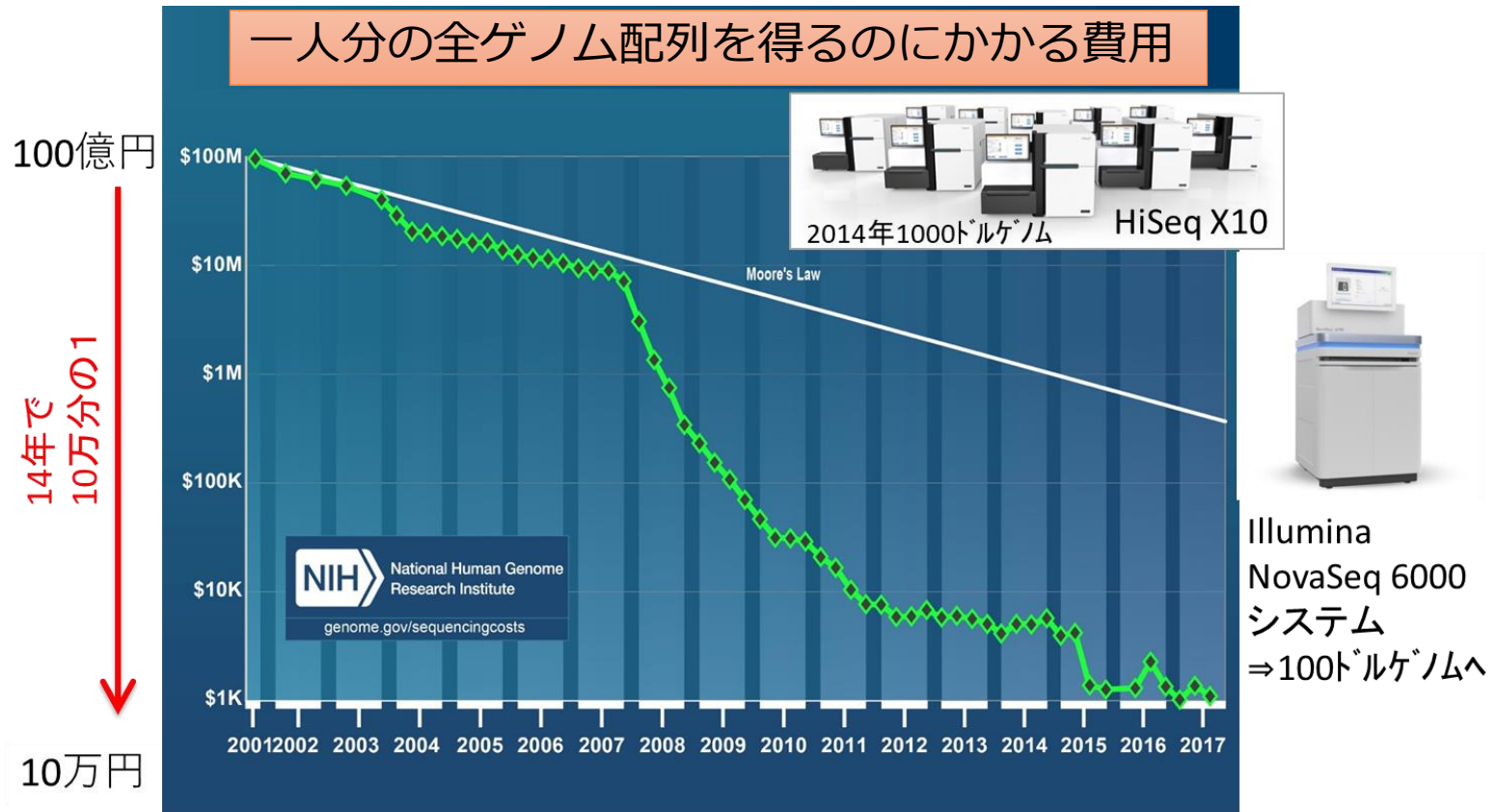


Integrative HMP (iHMP) Research Network Consortium,  
*Cell Host & Microbe*, (2014).

Price *et al.*, *Nat Biotechnol*, (2017)

# 様々な生体データが大量に取得できるようになった背景

計測装置の飛躍的性能向上・計測コスト低下



様々な生体データが大量に取得できるようになった背景には、DNAの塩基配列（シーケンス）などを計測する次世代シーケンサー(NGS)の飛躍的性能向上があります。一人分の全ゲノム情報を得るのにかかるコストは14年間で10万分の1になりました。RNAなどもNGSで計測できます。腸内細菌叢のゲノム（メタゲノム）なども読むことができます。これらのシーケンスデータからがん細胞特有の変異・変化を見つけることで、がんの原因や治療法に関する情報を得ることができます。

# シークエンサーって、 どんなデータ？

## 生のサンプル

DNAとして抽出  
ATCGの4種の文字



次世代シークエンサー



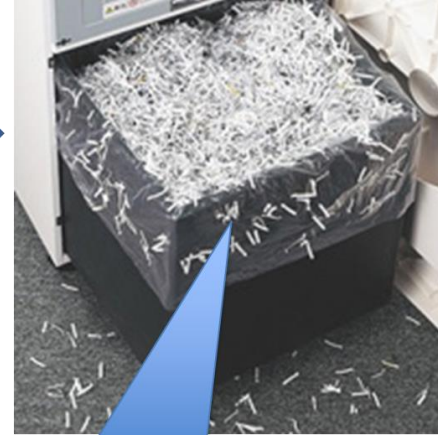
## 得られるデータ

ゲノムシュレッター

100文字ぐらいの断片になった

## 21億ピース

の文字列断片がコンピュータに  
吐き出される



ATCCGGTAAAT.....TTCA

100~150塩基

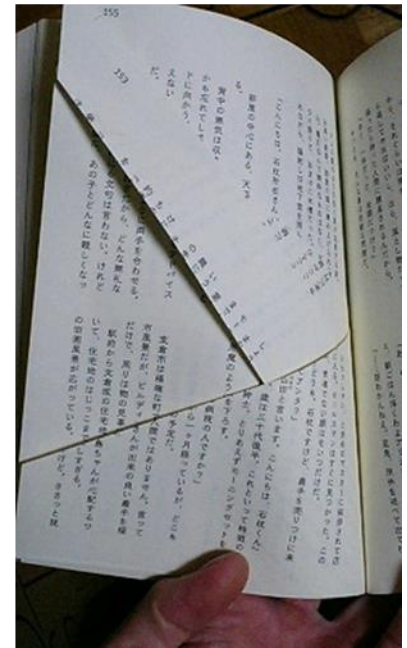
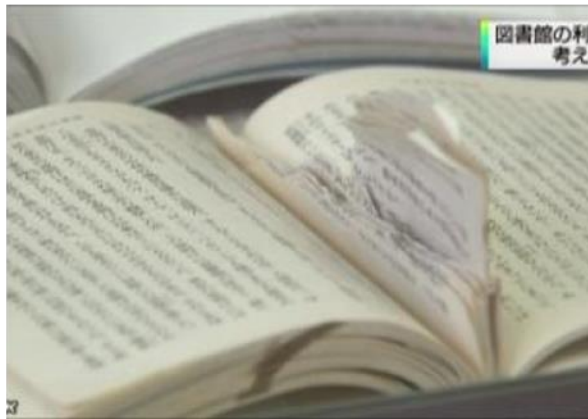
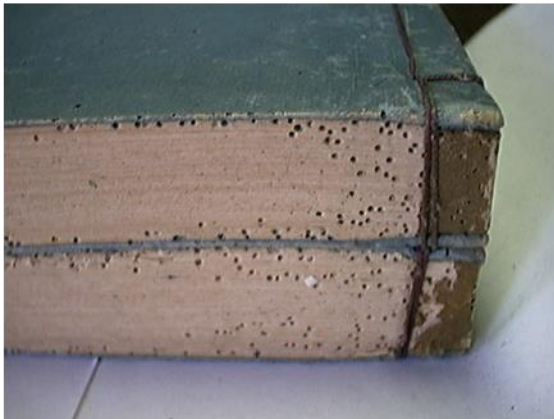
※全ゲノムシークエンスを想定

一人分のゲノムは、**4種類の塩基**（A、T、G、C）からなる塩基対、**約30億個分**からなります。現在主流のシークエンサーでは、一度にはDNA中の**約100塩基対分**しか読めません。これは**ATGCの4文字**からなる**30億文字分**の**文章が書かれた書類**の束が、シュレッターにかけられ**100文字ずつの断片**になって出てくるようなものです。通常、がん細胞特有の変異を見つける場合、**がん細胞由来の書類を40コピー分、正常細胞由来の書類を30コピー分**、シークエンサーで読み取ります。つまり**21億ピースの文字列断片**が得られます。ここから各ピースがゲノム上のどこにあったかを正確に決定し、**がん細胞と正常細胞の違いを見つけ出す（変異を検出する）**必要があります。

# 「変異」を探し出せ

1文字変異  
挿入や欠失  
転座、逆位・・・

本にたとえば





# 変異検出にチャレンジ

P53遺伝子(最も有名ながん抑制遺伝子)の一部。一文字だけ配列が異なります。

## 正常細胞のDNA配列

```
gatgggattg gggttttccc ctccatgtg ctcaagactg gcgctaaaag ttttgacctt
ctcaaaagtc tagagccacc gtccaggagg caggtagctg ctgggctccg gggacacttt
gcgttcgggc tgggagcgtg ctttccacga cggtagacag ctccctgga ttgggtaagc
tcctgactga acttgatgag tcctctctga gtcacgggct ctgggctccg tgtattttca
gctcgggaaa atcgctgggg ctgggggtgg ggcagtgggg acttagcgag ttgggggtg
agtgggatgg aagcttggct agagggatca tcataggagt tgcattgttg ggagacctgg
gtgtagatga tggggatggt aggccatcc gaactcaag ttgaacgcct aggcagagga
gtggagcttt ggggaacctt gagccggcct aaagcgtact tctttgcaca tccaccgggt
gctgggcgta gggaatccct gaataaaaag atgcacaaag cattgaggtc tgagactttt
ggatctcgaa acattgagaa ctcatagctg tatatttttag agcccatgpc atcctagtga
aaactggggc tccattccga aatgatcatt tgggggtgat cgggggagcc caagctgcta
aggtaaccaca acttccggac ctttgcctt cctggagcga tctttccagc gagccccgg
ctccgctaga tggagaaaat ccaattgaag gctgtcagtc gtggaagtga gaagtgttaa
accaggggtt tgcccgcag gccgaggagg accgtcgcaa tctgagaggc cggcagccc
tgttattgtt tggctccaca ttacatcttgc agcagcattt cgggttctt
tttgccggag cagctcacta ttcaccgat gagaggggag gagagagaga gaaaatgtcc
tttaggccgg ttcctcttac ttggcagagg gaggctgcta tctccgctga geatttctt
ttctggatta cttagttagt gcctttgcaa aggcaggggt atttgttttg atgcaaacct
caatccctcc ccttctttga atgggtgtgc ccaccccggc ggtgcctgc aacctaggcg
gacgctacca tggcgtgaga cagggagggga aagaagtgtg cagaaggcaa gcccgagggt
atthtcaaga atgagtatat ctcatcttcc cggaggaaaa aaaaaaagaa tgggtacgtc
tgagaatcaa atthtgaag agtgcattga tgggtcgttt gataatthgt cggaaaaaca
atctacctgt tatctagctt tgggctaggc cattccagtt ccagacgcag gctgaacgtc
gtgaagcggga aggggcgggc ccgcaggcgt ccgtgtggtc ctccgtgcag cctccggcc
cgagccgggt cttcctggtg ggaggcggaa ctccaattca tttctcccgc tgcccattt
cttagctcgc ggttgtttca ttcocgagtt tottccatg cacctgcgcc gtaccggcca
ctttgtcccg tacttaagtc atcttttcc taaatcgagg tggcaattac acacagccc
agtgcacaca gcaagtgcac aggaagatga gttttggccc ctaaccgctc cgtgatgctt
accaagtcaac agaccctttt catcgtccca gaaacgtttc atcaagcttc ttcccagctg
attcccagacc ccacctttat ttgatctcc ataaccatth tgctgtttgg agaacttcat
atagaatgga atcaggctgg gcgctgtggc tcaagcctgc actttgggag gccgaggcgg
cgggattact tgaggatagg agttccagac cagcgtggcc aacgtgggta atcccgtct
ctactaaaaa atacaaaaat tagctggggc tgggtgggtgc ctgtaatccc agctattcgg
gagggtgagg caggagaaat gcttgaaccc gggaggcaga ggttgcagtg agccaagatc
```

## がん細胞のDNA配列

```
gatgggattg gggttttccc ctccatgtg ctcaagactg gcgctaaaag ttttgacctt
ctcaaaagtc tagagccacc gtccaggagg caggtagctg ctgggctccg gggacacttt
gcgttcgggc tgggagcgtg ctttccacga cggtagacag ctccctgga ttgggtaagc
tcctgactga acttgatgag tcctctctga gtcacgggct ctgggctccg tgtattttca
gctcgggaaa atcgctgggg ctgggggtgg ggcagtgggg acttagcgag ttgggggtg
agtgggatgg aagcttggct agagggatca tcataggagt tgcattgttg ggagacctgg
gtgtagatga tggggatggt aggccatcc gaactcaag ttgaacgcct aggcagagga
gtggagcttt ggggaacctt gagccggcct aaagcgtact tctttgcaca tccaccgggt
gctgggcgta gggaatccct gaataaaaag atgcacaaag cattgaggtc tgagactttt
ggatctcgaa acattgagaa ctcatagctg tatatttttag agcccatgpc atcctagtga
aaactggggc tccattccga aatgatcatt tgggggtgat cgggggagcc caagctgcta
aggtaaccaca acttccggac ctttgcctt cctggagcga tctttccagc gagccccgg
ctccgctaga tggagaaaat ccaattgaag gctgtcagtc gtggaagtga gaagtgttaa
accaggggtt tgcccgcag gccgaggagg accgtcgcaa tctgagaggc cggcagccc
tgttattgtt tggctccaca ttacatcttgc agcagcattt cgggttctt
tttgccggag cagctcacta ttcaccgat gagaggggag gagagagaga gaaaatgtcc
tttaggccgg ttcctcttac ttggcagagg gaggctgcta tctccgctga geatttctt
ttctggatta cttagttagt gcctttgcaa aggcaggggt atttgttttg atgcaaacct
caatccctcc ccttctttga atgggtgtgc ccaccccggc ggtgcctgc aacctaggcg
gacgctacca tggcgtgaga cagggagggga aagaagtgtg cagaaggcaa gcccgagggt
atthtcaaga atgagtatat ctcatcttcc cggaggaaaa aaaaaaagaa tgggtacgtc
tgagaatcaa atthtgaag agtgcattga tgggtcgttt gataatthgt cggaaaaaca
atctacctgt tatctagctt tgggctaggc cattccagtt ccagacgcag gctgaacgtc
gtgaagcggga aggggcgggc ccgcaggcgt ccgtgtggtc ctccgtgcag cctccggcc
cgagccgggt cttcctggtg ggaggcggaa ctccaattca tttctcccgc tgcccattt
cttagctcgc ggttgtttca ttcocgagtt tottccatg cacctgcgcc gtaccggcca
ctttgtcccg tacttaagtc atcttttcc taaatcgagg tggcaattac acacagccc
agtgcacaca gcaagtgcac aggaagatga gttttggccc ctaaccgctc cgtgatgctt
accaagtcaac agaccctttt catcgtccca gaaacgtttc atcaagcttc ttcccagctg
attcccagacc ccacctttat ttgatctcc ataaccatth tgctgtttgg agaacttcat
atagaatgga atcaggctgg gcgctgtggc tcaagcctgc actttgggag gccgaggcgg
cgggattact tgaggatagg agttccagac cagcgtggcc aacgtgggta atcccgtct
ctactaaaaa atacaaaaat tagctggggc tgggtgggtgc ctgtaatccc agctattcgg
gagggtgagg caggagaaat gcttgaaccc gggaggcaga ggttgcagtg agccaagatc
```





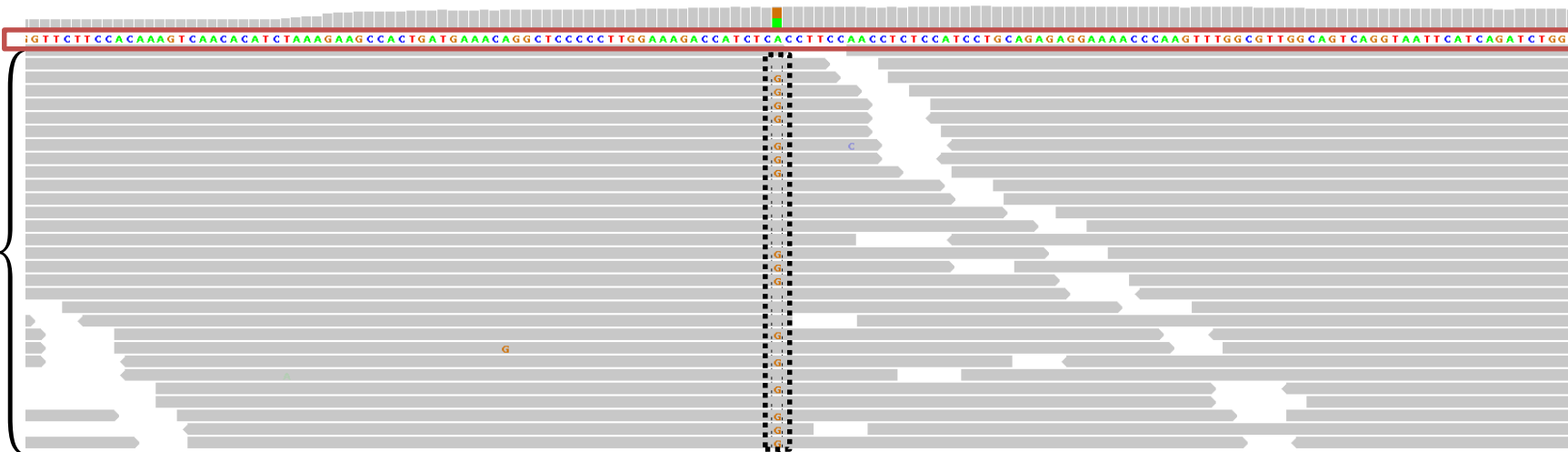
# DNA文字列ピースデータからの変異検出原理

文字列ピースデータ  
(約100文字; 21億ピース)



ステップ1:  
スーパーコンピュータを使って  
ヒトゲノム標準配列(30億文字)  
上でマッチする場所を探索します  
(アライメント)。

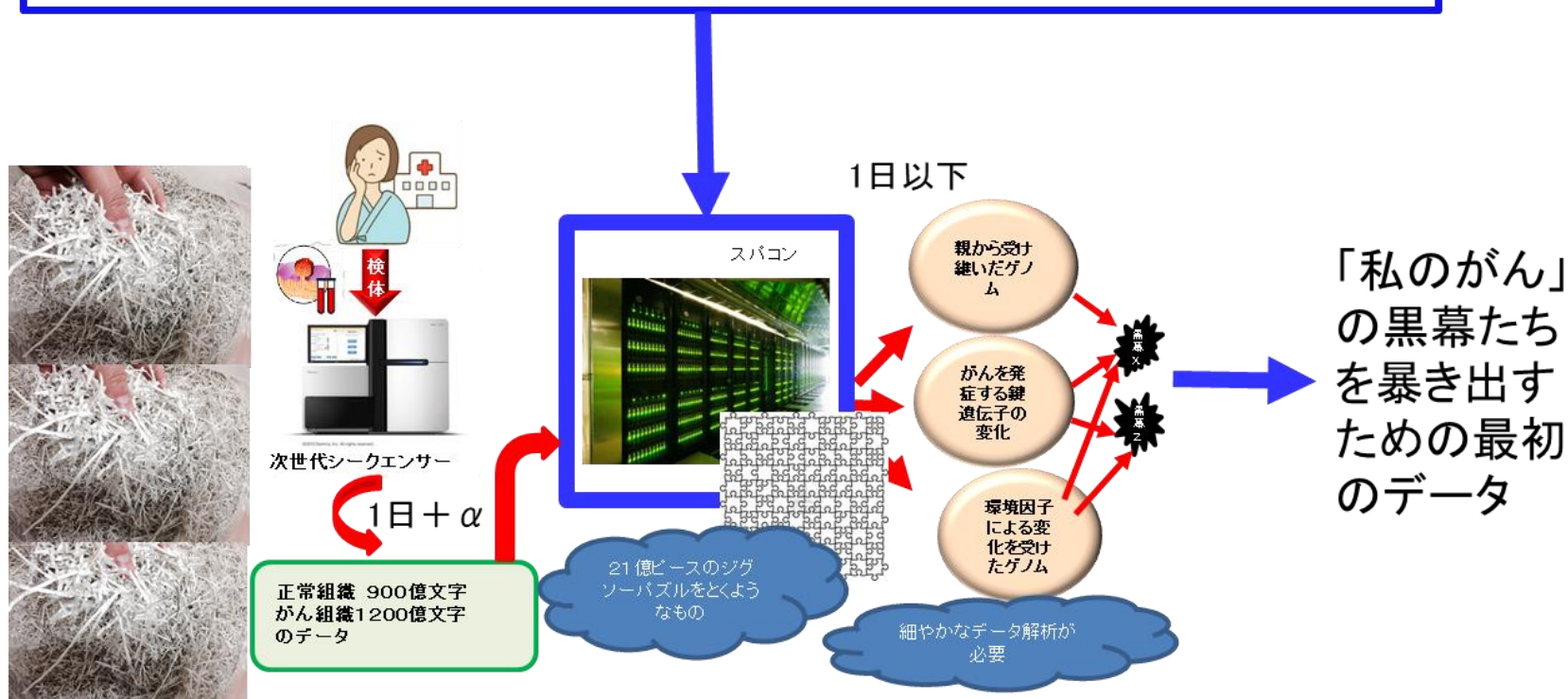
ヒトゲノム標準配列(30億文字)



アライメントされた文字列ピースデータ。  
上段の標準配列と同じ文字(塩基)で  
あれば灰色で表示。

ステップ2:  
アライメントされた文字列ピースデータ中で、  
標準配列と異なる文字が複数回観測されてい  
る場所を検出します。この例では標準配列が  
Aのところ、Gに置き換わっています。  
(一塩基変異)

スパコンで21億ピースのジグソーパズルを解き、がんのシステム異常の原因を暴き出さねばならない！



一塩基変異だけでなく、欠失、増幅、転座などの変異もあります。それらをスーパーコンピュータと変異検出アルゴリズム（多くは統計的手法に基づく）を用いて、高精度かつ高速に検出する必要があります。最近ではGPUと呼ばれる演算装置を活用した計算の高速化も進んでいます。

# 変異検出の難しい点（1）

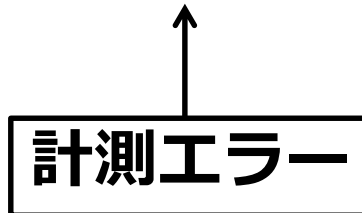
## 計測エラーの混入

- DNA配列のシーケンサーでの計測時に**計測エラー**が入ることがあります
  - 本当の変異がある位置を検出する際のノイズとなります

**ATCGGACCATGTCCAATCA** 本当のDNA配列



**ATCGGACCATGTCCA****G****TCA** 計測されたDNA配列  
(エラー有)

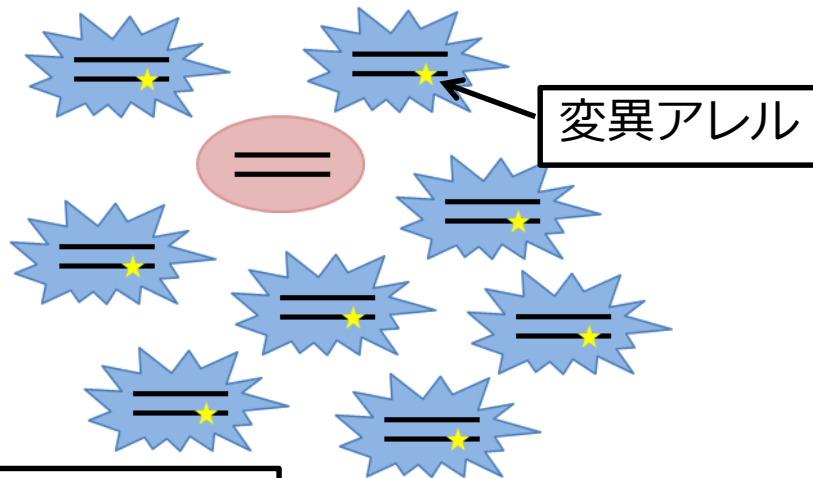


# 変異検出の難しい点 (2)

## がん細胞含有率が低い場合

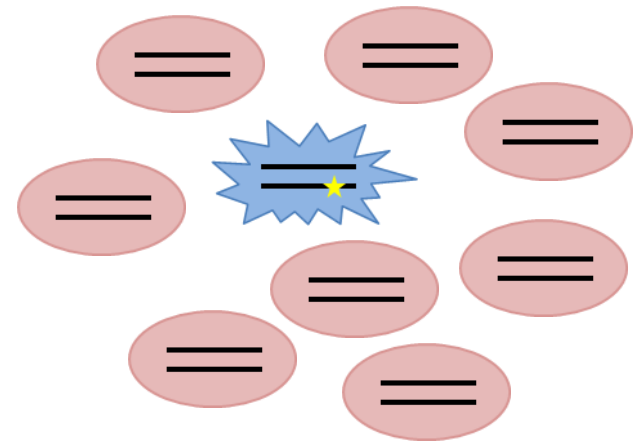
- サンプル中のがん細胞含有量が少ない場合があります
  - がん組織中には正常な細胞も含まれます
  - **変異アレル観測割合**が相対的に少なくなり、エラーとの区別が難しくなります

高がん細胞含有サンプル

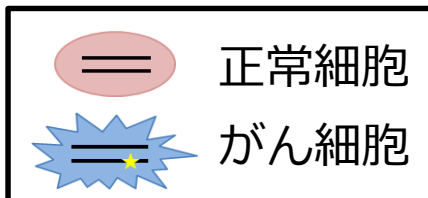


変異アレル(★)割合 : 44%

低がん細胞含有サンプル



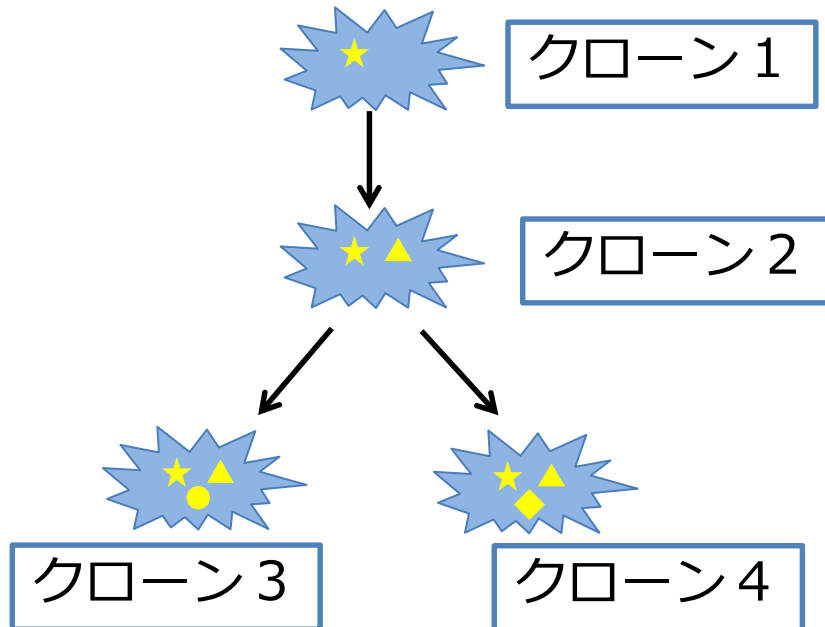
変異アレル(★)割合 : 5.6%



# 変異検出の難しい点 (3)

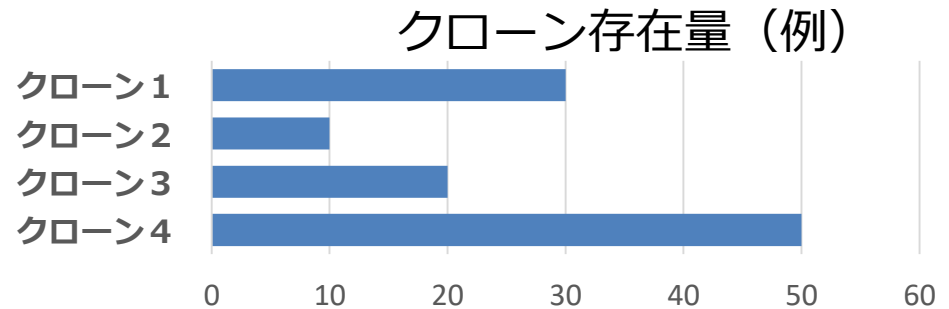
## がん細胞クローンの存在

- がん細胞は**進化**します。新しい変異を獲得したがん細胞クローンが出現することがあります
  - 変異の獲得時期、クローンの存在量に応じて変異の観測頻度が変わります
  - 存在量の少ないクローン固有の変異の検出は難しい

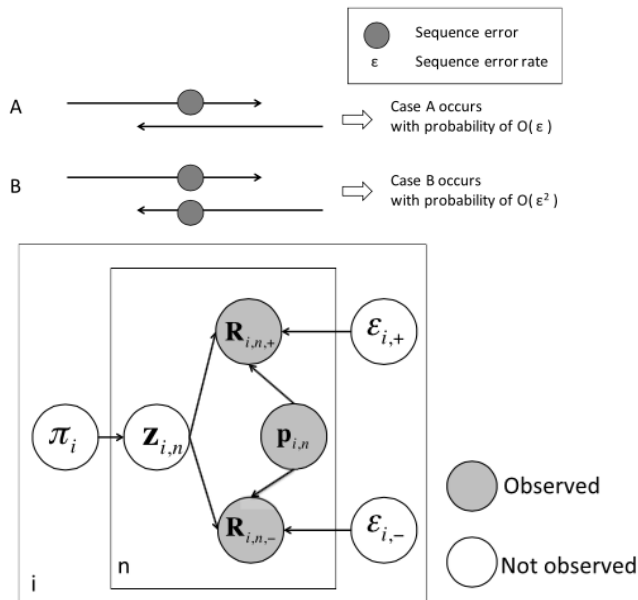


変異

|       | ★ | ▲ | ● | ◆ |
|-------|---|---|---|---|
| クローン1 | ✓ |   |   |   |
| クローン2 | ✓ | ✓ |   |   |
| クローン3 | ✓ | ✓ | ✓ |   |
| クローン4 | ✓ | ✓ |   | ✓ |



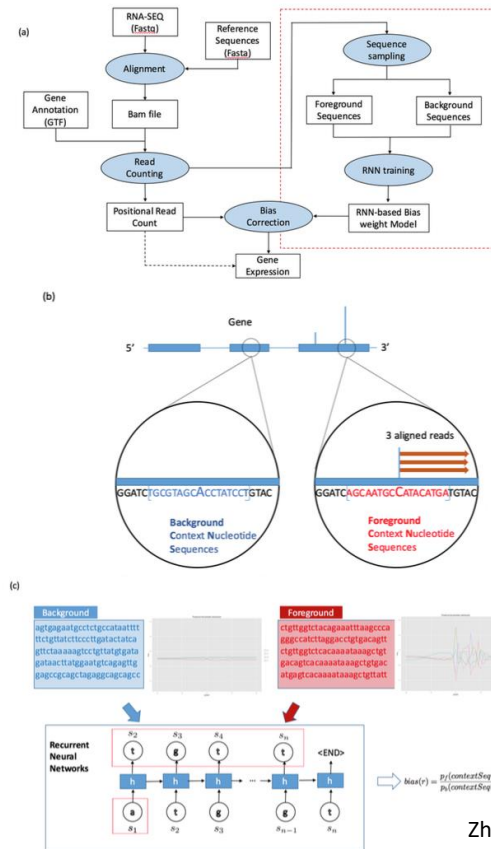
# ベイズ統計モデル



$$\begin{aligned}
 & p(\mathcal{R}_i, \mathcal{Z}_i | \gamma_i, \alpha_{i,+}, \alpha_{i,-}) \\
 &= p(\pi_i | \gamma_i) p(\epsilon_{i,+} | \alpha_{i,+}) p(\epsilon_{i,-} | \alpha_{i,-}) \\
 & \quad \cdot \prod_n p(\mathbf{R}_{i,n,+}, \mathbf{R}_{i,n,-} | \mathbf{z}_{i,n}, \epsilon_{\pm,i}, \pi_i, \mathbf{p}_{i,n}) p(\mathbf{z}_{i,n} | \pi_i) \\
 &= p(\pi_i | \gamma_i) p(\epsilon_{i,+} | \alpha_{i,+}) p(\epsilon_{i,-} | \alpha_{i,-}) \\
 & \quad \cdot p(\mathbf{R}_{i,+}, \mathbf{Z}_{i,+} | \epsilon_{i,+}, \pi_i) \cdot p(\mathbf{R}_{i,-}, \mathbf{Z}_{i,-} | \epsilon_{i,-}, \pi_i) \\
 & \quad \cdot p(\mathbf{R}_{i,\pm}, \mathbf{Z}_{i,\pm} | \epsilon_{i,\pm}, \pi_i)
 \end{aligned}$$

Moriyama et al., IEEE Trans Nanobio, 2017  
 Hayashi et al., BMC Genomics, 2018  
 Hayashi et al., J Comput Biol, 2019  
 Moriyama et al., Bioinformatics, 2019  
 Moriyama et al., LNCS, (in press)

# 深層学習モデル



Zhang et al., BMC Genomics, 2017  
 Zhang et al., IEEE BIBM, 2017  
 Konishi et al., BMC Bioinfo, 2019  
 Zhang et al., BMC Bioinfo, (in press)

前のページで説明した困難を克服して、高精度に変異を検出する必要があります。そのため、データの生成機構をモデル化したベイズ統計に基づく変異検出手法や、深層ニューラルネットワークに基づく変異検出手法などの開発研究を行っています。



# おわりに

- がんゲノムの変異を検出するため用いる、観測データ、変異検出の原理、計算資源、困難な点、それを克服するためのデータ解析手法などについて概説しました
- 次世代シーケンサーによる観測データは、ゲノム変異の検出に加えて、個人の免疫遺伝子型決定や、免疫細胞クローンの多様性の推定にも用いられ、がん免疫療法の開発にとっても不可欠なものとなっています
- 今後もシーケンサーの性能向上に伴い、新たな種類のデータが産生されると期待されます。それに応じて、新たな解析手法の開発を行う必要があります
- また大量のデータ解析結果を、いかに医療に有用な情報へ翻訳し還元するかということも課題となっており、人工知能を活用する研究が進みつつあります
- それに関しては、また別の機会にご紹介できればと思います