

第2回 SPARC Japan セミナー2018

「オープンサイエンス時代のクオリティコントロールを見通す」

学術コミュニケーションのエコシステムの今後 ～arXivの現状から考える～

武田 英明

(国立情報学研究所)

講演要旨



インターネットの利用とオープンアクセス活動の浸透により、学術コミュニケーションのエコシステムが大きく変わってきている。伝統的なジャーナル投稿掲載という方法以外にも様々な方法が生まれている。現在の学術コミュニケーションのエコシステムを概観し、その中での preprint server の立ち位置を確認する。特に代表的 preprint server である arXiv を取り上げ、現状と課題を確認する。



武田 英明

SPARC Japan運営委員会委員長。

http://www.nii.ac.jp/faculty/informatics/takeda_hideaki/

本日はいつもの話と少し違って arXiv、プレプリントを中心に話をしたいと思います。なぜ私がプレプリントの話をするかというと、10月2日に arXiv の Member Advisory Board (MAB) というミーティングに

行ってきたからです(図1)。arXiv というのはメンバー制度になっていて、大学等がメンバーです。日本の場合は日本コンソーシアムをつくっていて、幾つかの大学が入っています。そのコンソーシアムの代表が毎回 MAB に出ています。コンソーシアムの代表は、本当は京都大学の引原隆士先生なのですが、代表が代わる時期で、コンソーシアム代表代理という形で私が MAB に参加してきました。

MAB は1年に1回行われます。Member Advisory Board と、もう一つ Scientific Advisory Board (SAB) というボードがあり、この二つで arXiv はマネジメントされています。MAB の参加者の多くは図書館で、一部は Jisc などの研究助成機関です。研究者代表の方もいるのですが、今回は欠席でした。場所はコーネル大

arXivのMember Advisory Board (MAB) Meetingに行ってきました。



- 10/2 @Olin Library, Cornell University
- この会議はarXivのメンバーのうち、コンソーシアムなどの代表者などを招集して、1年に1度開催。
- 参加者
 - 多くが図書館：
 - Los Alamos Library, TIB, U of California, Ohio State Univ, U of Amsterdam, U of Queensland, U of Sydney
 - 研究助成機関
 - Jisc, Simons Foundation
 - 研究者代表: professor at KTH (欠席)



(図1)

学です。私は初めてコーネル大学に行きましたが、きれいで広大なキャンパスでした。

arXivの現状と課題

現在、arXivには約140万論文があり、1日に約600件、新しい投稿が来ます。1秒に7論文がダウンロードされます。モデレーターという、論文のチェックをする人が162人います。プレプリントは論文査読はしません。ただし、分野ごとにモデレーターのグループがあって、そのモデレーターが、その分野の論文として適切かどうかを評価します。それに通れば1~2日で掲載されるようになっています。

最近の傾向を見ると、2012年のころはサブミッションが年間84,000件でしたが、2017年は約120,000件で、5割増ぐらいのペースで増えていることがわかります(図2)。ポイントは、この増えた量がどこから来ているかというところです。

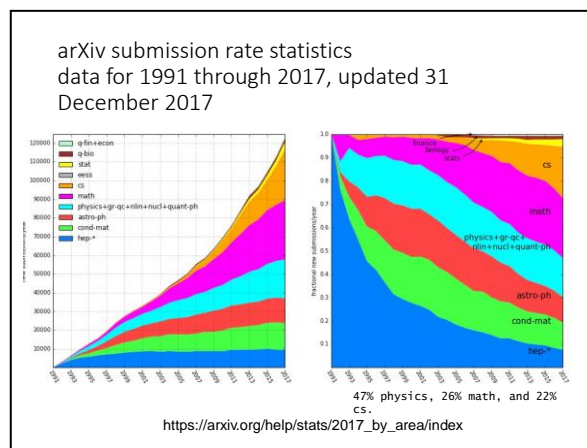
arXivはできて30年たっていますが、図3の統計は1991年から始まっています。下から、青色が高エネルギー物理学、緑色が固体物理学、赤色が天体物理学です。これに水色を加えた下から四つが物理系です。さらに、紫が数学、オレンジがコンピュータサイエンスです。小さく見えるのが最近入ってきた統計など、広い意味での経済系です。

これを見て分かるように、元々の始まりである高エネルギー物理学はほとんど変わらない、減りもしなければ大して増えもしない状況です。それ以外の物理は

やや増えています。数学も比較的増えています。近年増加率が大きいのは、何といたってもコンピュータサイエンスです。2017年は投稿の22%がコンピュータサイエンスとなり、コンピュータサイエンスが顕著に増加しています。高エネルギー物理学は、相対的にはarXivの中のウエートとしては顕著に下がってきているということがわかります。とはいえ、物理全体で、5割弱ぐらいは占めているという状況にあります。

これが分野としての分析です。これは図の下に記載されたサイトで手に入ります。

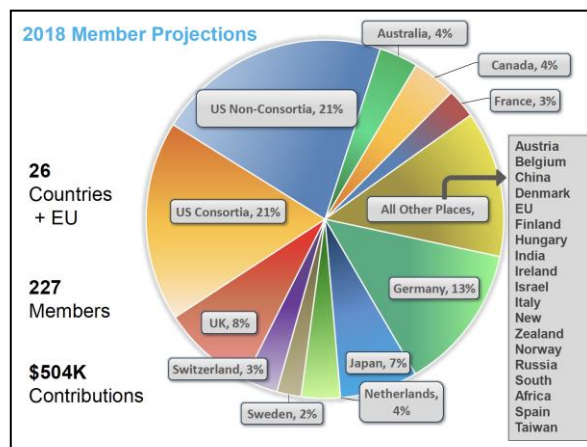
図4は、arXivのメンバーはどここの国から来ているかという分布です。US Consortiaが、アメリカの大学等がコンソーシアムを組んで加盟しているメンバーです。US Non-Consortiaがそれ以外の独立で入っているメンバーです。それらを合わせてアメリカで約40%です。そしてイギリスなど主要なヨーロッパの国があ



(図3)



(図2)



(図4)

り、アメリカとヨーロッパ以外には、中国やオーストラリアがあります。日本も7%入っていて、US、ヨーロッパ以外では最大となっています。

現在、200強のメンバーが加盟しています。このメンバーはお金を払います。arXivは多く投稿されている機関にメンバーになることを依頼しています。東京大学に所属する人がarXivにたくさん投稿していれば、東京大学にメンバーになってほしいと依頼します。ただし、arXivのメンバーになるのは強制ではありません。arXivはご存じのように誰でも、所属があろうとなかろうと投稿できます。ただ、財政のためにこのような仕組みを取っています。

どこから投稿されているかという、メンバーの国と比例するように、アメリカが約40%でヨーロッパは約30%です(図5)。中国はメンバー数に比べて投稿数が多く、3番目となっています。これは課題です。

運営体制は、大きく、執行チームと運営チームと技術アドバイザーとなっています(図6)。その他にアドバイザーのような形で、Paul Ginspargがいます。彼はarXivをつくった人で、今、コーネル大学のコンピュータサイエンス学科にいます。

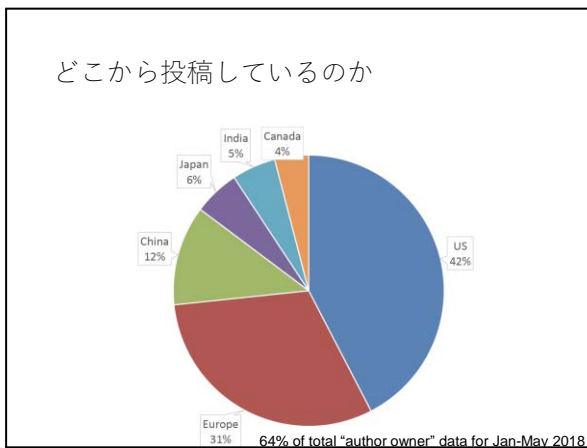
執行チームのOya Riegerさんはライブラリアンです。今年まではarXivはコーネル大学図書館が運営するという立場を取っています。だから、執行チームのトップのOyaさんはライブラリアンなのです。

既にカレントアウェアネスでも報告が出たと思いますが、来年からこの体制が変わって、コンピュータサ

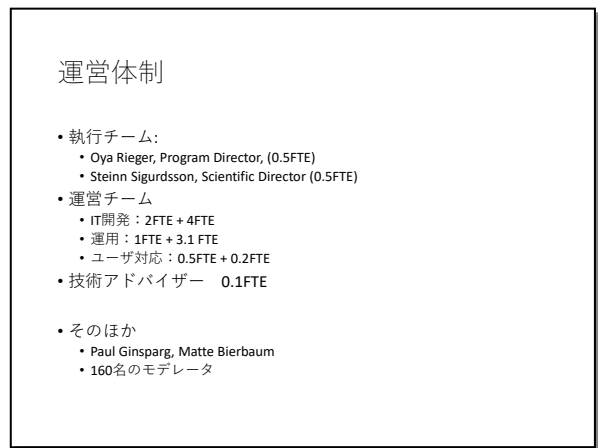
イエンス学科に移ります。われわれも心配したのですが、この運営体制そのものはほとんど変わらなくて、ただ、コーネル大学での責任を持つ部局が変わるということです。チームそのものが今、図書館の中にあるのですが、当面は物理的にも図書館の中にあるそうです。そういう意味では、ユーザの立場からだと、すぐ来年1月から何か変わるということはありません。

今まで、開発のことは図書館が責任を持つという体制で、執行チームと図書館が並行にあって、複雑だったのですが、コンピュータサイエンスに変わったら、今度は執行チームが開発チームを管理するという、むしろシンプルになると説明されました。Paul Ginspargのアドバイスを受けて、よりITシステムとしての改善を図りたいということが一つのモチベーションとなったそうです。

ビジネスモデルは重要なところで、これはサービスである以上、お金がないと回りません(図7)。一体どうやって回しているか、大きく言うと三つのリソースからできていて、まずはコーネル大学自体資金を出しています。現金も出しているし、Oyaさんのようなコーネル大学所属の人に、人件費という形で出しています。それから、いわゆる財団系のもの、Simons Foundation等が資金を日本円にすると4,000万円ぐらい出しています。あとは先ほど言ったメンバーから集めたお金が5,000万円ぐらいです。3分の1ずつを負担しているようなイメージで動いているのが現在の仕組みです。総予算が約1億5,000万円です。



(図5)



(図6)

この金額が妥当かどうかというところは結構な問題なものです。arXivを使ったことがある方は分かると思いますが、arXivは実は結構、古くさいシステムです。それを今、ニュージェネレーションシステムに非常に時間をかけて移行しようとしています。歴史があるので、とても複雑なシステムになっています。LaTeXから入って、そのままPDFが生成されるので便利なのですが、その分複雑なシステムになっていて、それがシステムをリニューアルするのにネックになっています。

今の話をまとめます。arXivの分野としてコンピュータサイエンスが正式に入ったのが2012年ぐらいで、まだ5年ほどしかたっていません。それからさらに最近、統計や一部の経済系も入ってきています。先ほど表で示したように、arXivは投稿数も増えているので、プレプリントサーバーを維持して、分野を順次増やしているという面においてはarXivは順調に成長しています。

また、arXivがつくったプレプリントの文化およびビジネスモデルが、OA時代を迎えて、分野を超えて高く評価されているということは大きいです。

図8が最近の動きです。arXivは1991年に誕生しています。そこから20年以上、大した動きはなく、2013年ごろから急激にプレプリントサーバーが他の分野で増えてきました。PeerJのプレプリントが出てきたのが2013年で、2016年以降、立て続けに出てきて、2016年にbioRxiv、engrXiv、SocArXiv、PsyArXiv

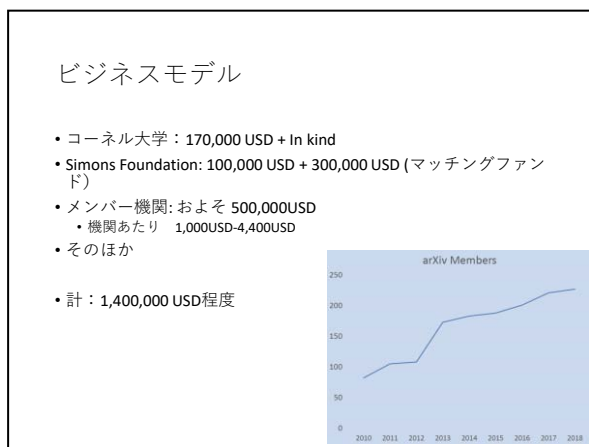
が出てきています。2017年にはLawArXiv、AgriXiv、ChemRxivなども出てきています。

このうち、後ろにxivと付いているものは、基本的にはarXivをリスペクトした、あるいはインスパイアされたとそれぞれのサイトのページに書いてあります。arXivのモデルをわれわれの分野でも使いたいと考えて立ち上げたと書いてあります。特に、赤いものがarXivにインスパイアされたものです。

ここに違うもの、SSRNがあります。これは約2年前にエルゼビアに買収されています。SSRNは元々、ソーシャルサイエンスのプレプリントサーバーでしたが、エルゼビアはこれをいろいろな意味で活用しようとしています。

その系列に当たるものがBioRNです。RN系列のバイオ版をエルゼビアが立ち上げたということのようです。

また、同じarXiv系でも後ろにOSFと書いてあるのが、これはOpen Science Frameworkで、Center for Open Scienceがつくっているプラットフォームに乗って、それを利用してサービスを提供しているプレプリントです。この五つは、マネジメントが違うだけでプラットフォームは同じものを使っています。それとは全く違う、インスパイアされたけれど、独自運営されているのは、bioRxiv、ChemRxivです。これはアメリカとドイツとイギリスの化学学会が共同運営をしています。面白いのは、プラットフォームはfigshareが提供するものを使っていると言っていることです。



(図7)



(図8)

2013年、特に2016年以降にたくさんのプレプリントが生まれていて、ただ、多くはarXivのつくった文化をむしろ自分の分野に取り入れたいという形で使われています。そういう意味で、arXivが30年間行ってきた活動は大変高く評価されています。

ところが、この30年で学術コミュニケーションが変わってきました。オープンアクセスの動きが強くなってきましたが、30年前はそのような動きはなかったのです。プレプリントサーバーの役割は変わってきています。arXivが変わった環境の中できちんと対応できるかということが、arXivという組織の課題となっています。

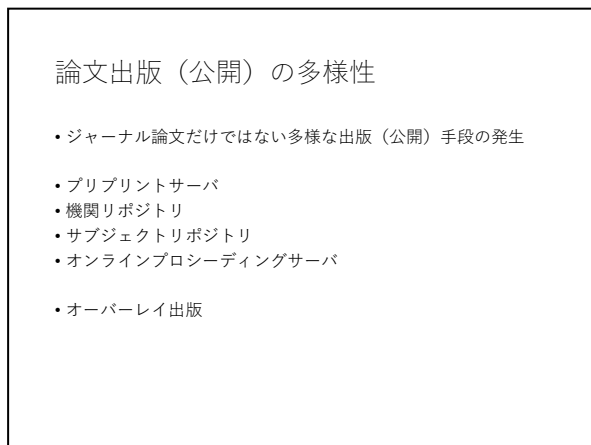
プレプリントサーバーの新しい役割

arXivができた当初、高エネルギー物理学では、新しい知見が発見されると、まずは関係する研究機関で情報を共有するという文化がありました。当初はファクスで送ったりするということがあったそうです。高エネルギー物理学の場合は、やっている研究所と大学が限られていて、世界中に拠点は数えられるほどなので、そういうことをやっていました。しかし、Paul Ginspargが、それでは効率が上がらないから、一つサーバーを置いて、そのレポートを全部蓄積すれば楽だろうと考えて始めたのがarXivです。大体の場合、そのレポートはいずれジャーナルにパブリッシュされるというものです。始まったころのプレプリントサーバーの役割は、比較的閉じたもので、あくまでも同業

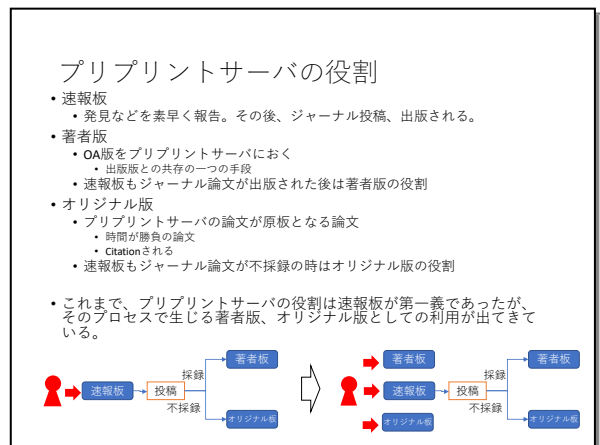
者が知らせ合うというサービスでした。

ところが、2017年、2018年においては、いろいろな論文の発表の仕方があって、その一部としてプレプリントサーバーが位置付けられるようになりました。これは、当初のプレプリントサーバーの役割と違う役割を持たせるようになったということです。だから、プレプリントサーバー自体の機能もそのようにならなければいけないというのが、今の課題です。具体的には、他のpublishing（出版・公開）の方法との役割分担、あるいは他のpublishingとの連携、さらには他のオープンアクセスへの積極的関与ということが求められています。

多様なpublishingのやり方が生まれてきました（図9）。では、具体的に何が違って来たかということですが、元々のプレプリントサーバーの役割というのは、まず速報版を出して、それがジャーナルに投稿されるというものです。採録されても消されることはないの、それは著者版として残ります（図10）。採録されないと、結果的にオリジナル版になります。つまり、ジャーナルには公開されないけれど、arXivには残るのです。これが元々のarXivの運用方法でした。それが逆に、初めからもう投稿しないつもりで出す人が出てきたのです。このプロセスに関係なく著者版を置きたいというだけのモチベーションでarXivを使う人も出てきました。つまり、速報版、著者版、オリジナル版が同じarXivの中に共存しているということになります。



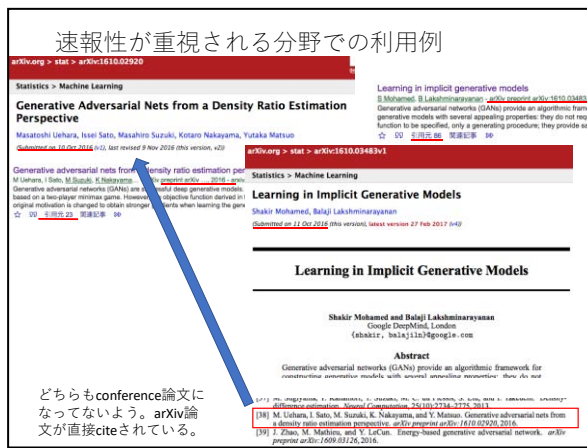
(図 9)



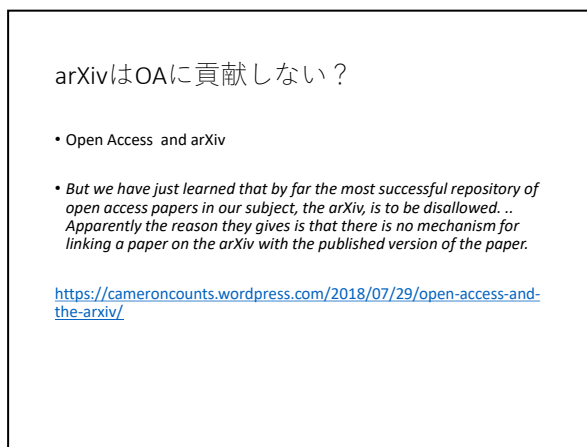
(図 10)

特に arXiv は速報性が重視される分野で使われるようになってきました。図 11 は昔、北本朝展さんが紹介したものが面白かったのもう一回引っ張ってきました。この論文が 2016 年 10 月 10 日に arXiv に載せられると、別の arXiv の論文が引用するのです。その論文の出版日が 2016 年 10 月 11 日ということで、前日の投稿を引用しているのです。

調べたところ、この二つの論文ともジャーナルには採録されていないようなのです。でも、そんなことは関係ないのです。Google Scholar でどれだけ使われているかを見てみると、例えば右の論文は 86 件引用されています。左の論文も 23 件引用されています。ですから、別にジャーナルに採録されるかどうかは関係なく、arXiv にあるだけである種の目的は達成しているのです。



(図 11)



(図 12)

arXiv は OA に貢献しない？

もう一つは、ジャーナル、カンファレンス連携ということがあります。これは最終版とのリンクがあること、投稿時に連携することです。

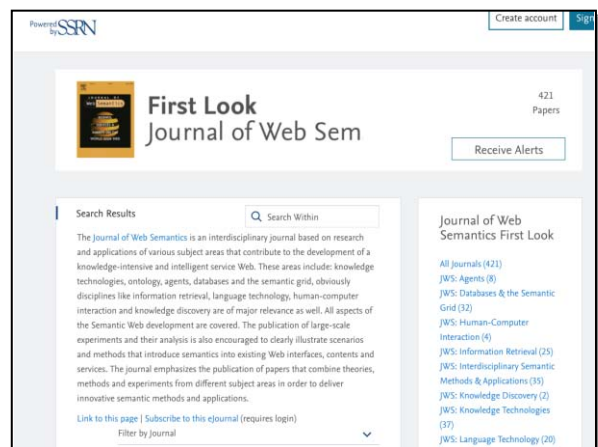
ある人が、arXiv は出版版とのリンクがないので、オープンアクセスの手段としてはカウントされないということを言っています (図 12)。実は、arXiv には明示的に出版版とリンクする仕組みがないのです。また、DOI もないのです。ここが問題ではないかということです。

後で坊農さんが説明する bioRxiv はきちんと DOI が付きます (図 13)。また、bioRxiv においてジャーナルに投稿するという仕組みが確立されています。従って、当然、投稿されると自動的に出版版とリンクが付くという仕組みを用意しています。

図 14 は私がエディターをしている『Journal of Web



(図 13)



(図 14)

Semantics』というエルゼビアが出しているジャーナルです。これは最近、投稿時は公開しないのだけでも、アクセプト時には著者版が SSRN に自動的に置かれるという機構を今度から導入すると言っています。今までのジャーナルはプライベートなサーバーに著者版を置いていたのですが、そうするとアクセスしづらいため、今までの 400 件の論文を全部含めて SSRN に置き換えたのです。

『Journal of Web Semantics』は、全くソーシャルサイエンスではなくてコンピュータサイエンスなのですが、なぜか SSRN に置くということです。これは恐らくエルゼビアがプレプリントサーバーを戦略的に使うための一つのトライアルです。こういうことが起きていて、これも著者版とのリンクです。

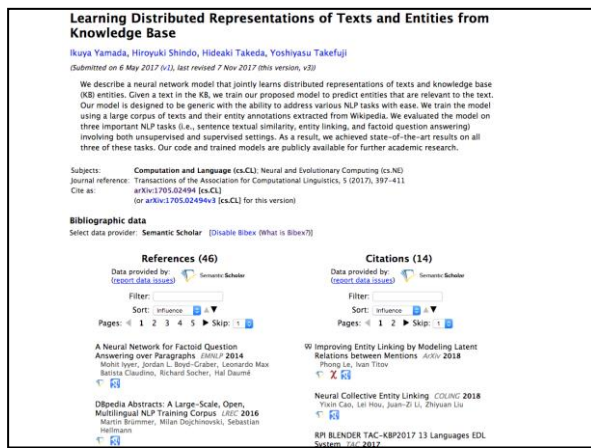
では、arXiv は今どうしているかという、DOI 付与、メタデータの再設計に関してもまだめどが立って

いません。それはテクニカルな理由です。システムが直せないというのです。メタデータがきちんとできないと言っていて、まだ見込みが立たないのです。

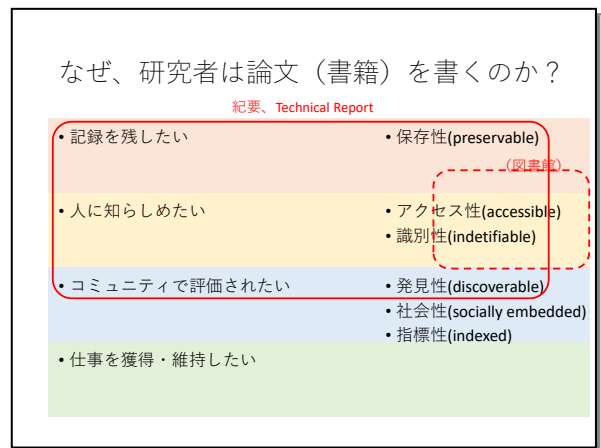
ジャーナル、カンファレンス連携に関しては、実験をしています。ACM と連携実験をしていて、まさに ACM のカンファレンスに投稿するときに、同時に arXiv に載せるという仕組みを今テスト中です。

今回の MAB にも ACM の人が来ていて、ACM および学会が arXiv とぜひ協力したいと。つまり、われわれはむしろ一緒に共存することが学会の使命だと思っているとはっきり言いました。ただ、arXiv 側がなかなかシステムの運用がままならないのが現状ということです。

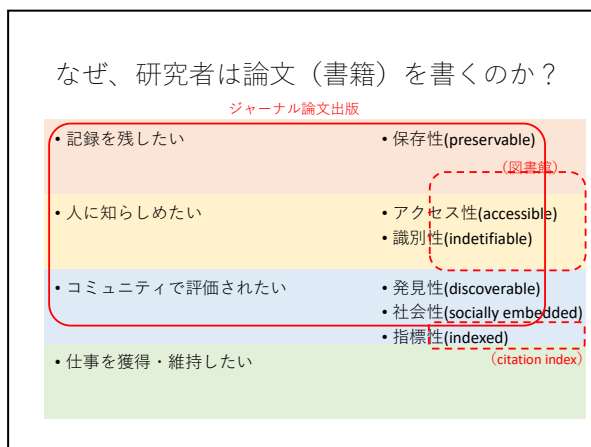
また、arXiv Labs というところで、少しだけシステム改善で、サイテーションが下に出るような仕組みを実験的に導入しています。皆さんもぜひ試してください。



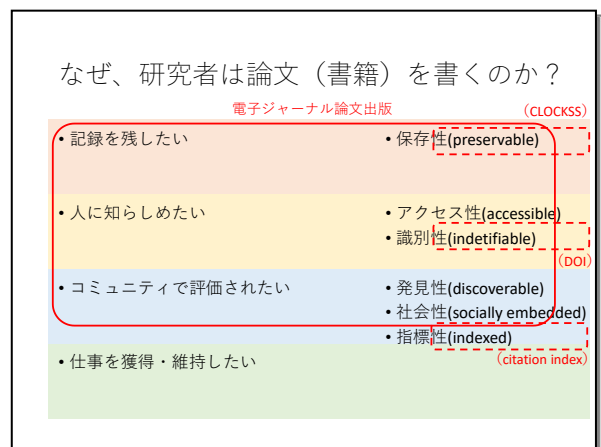
(図 15)



(図 17)



(図 16)



(図 18)

リファレンスとサイテーションが出るような仕組みができています (図 15)。

結論

なぜ研究者は論文を書くのか。それは、記録を残したい、人に知らしめたい、コミュニティで評価されたい、仕事を獲得したい・維持したい、これに尽きるのです。

17世紀は、それを往復書簡で行っていました。それを書籍に書くと、人に知らせることができます。でも、評価されるのはなかなか難しいので、書評などがあります。

ジャーナルというのはある種の素晴らしいシステムで、記録、人に知らしめる、コミュニティに評価されたいというのをワンセットにして、非常によく回りました (図 16)。紀要は図 17 のようになります。電子ジャーナルは役割は紙版と同じですが、DOIなどで、より識別性が良くなりました (図 18)。

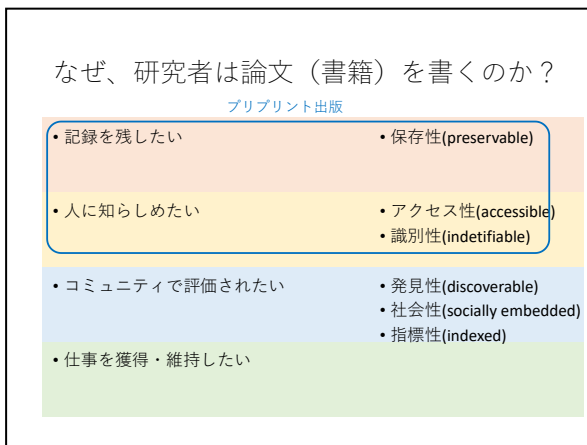
では、プレプリントではどうかというと、コミュニティで評価されたいというところまで入ってきません (図 19)。

そのときに、青い部分と赤い部分で役割分担をしようというのが今日の午前中のお話でもあり、プレプリントの新しい役割です (図 20)。評価されるといった部分は別の機能に持たせればいいのか。そういうことが今起きているのがジャーナル連携です。

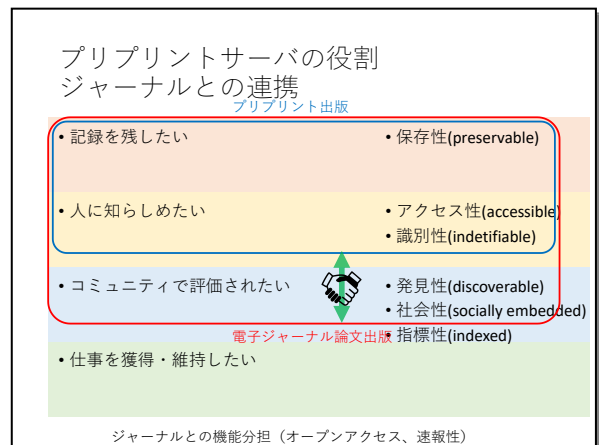
一方、直接オリジナル版でいいではないかという人

にとっては、むしろプレプリントサーバーの機能を強化してほしいというところがあります。機能を強化するということは、いろいろな指標性、アクセス性、識別性をプレプリントサーバー自体に持ってほしいという要求ということになります。

まとめると、プレプリントサーバーは求められる機能や役割が増えてきています。今回の MAB に出席して受けた印象は、arXiv はシステムはあまりモダンではなく、またマネジメントも比較的ゆっくりしているということです。ですので、現在の変化のスピードについていけるかという点を危惧しました。arXiv はサービス自身もサービスモデルも非常に期待されている一方、なかなか困難な段階になっていると感じました。



(図 19)



(図 20)