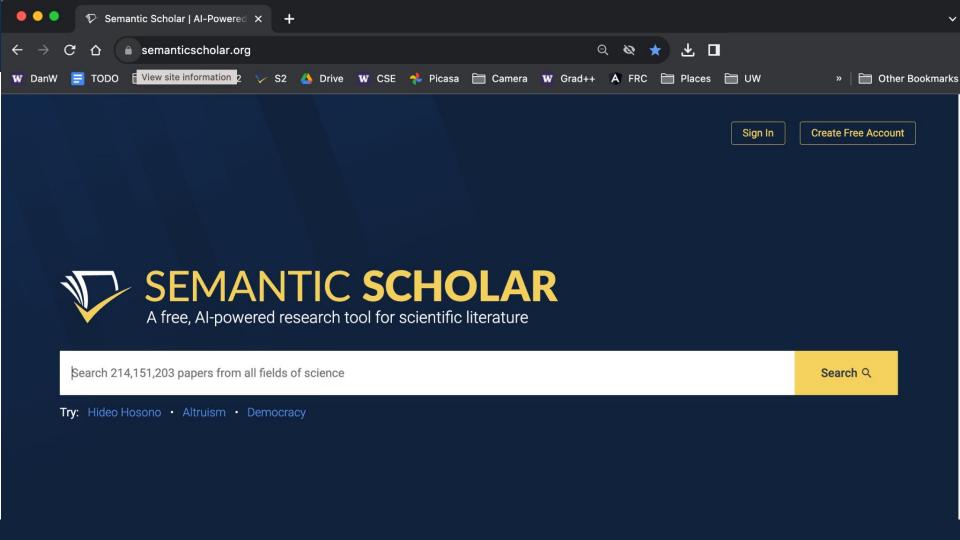
Al for Scientific Assessment

Daniel S. Weld

Allen Institute for Artificial Intelligence
University of Washington



Semantic Scholar Mission:

Accelerate Science Breakthroughs Using Al

Ai2 is a non-profit & Semantic Scholar is a free service

8M monthly active users

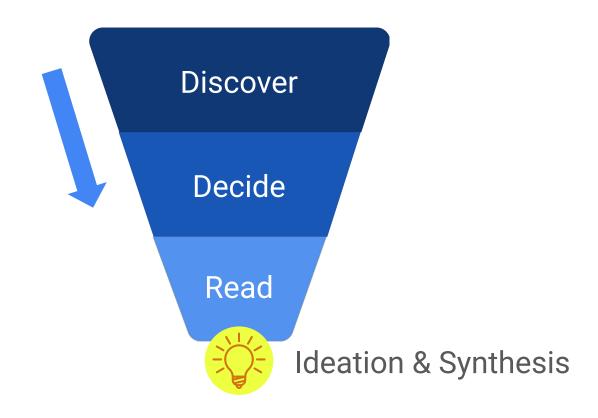
221M scientific paper index

30+ data sources

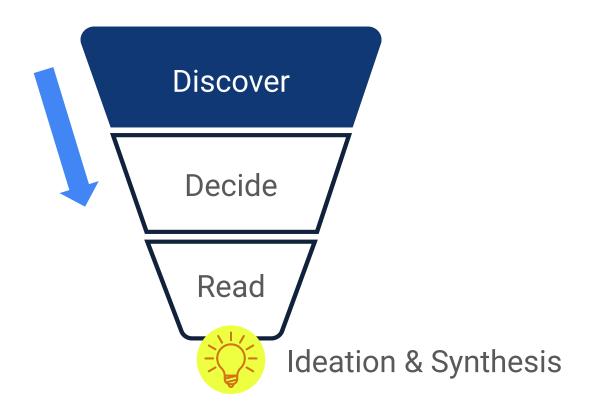
25+ ML models



The Scholarly Funnel

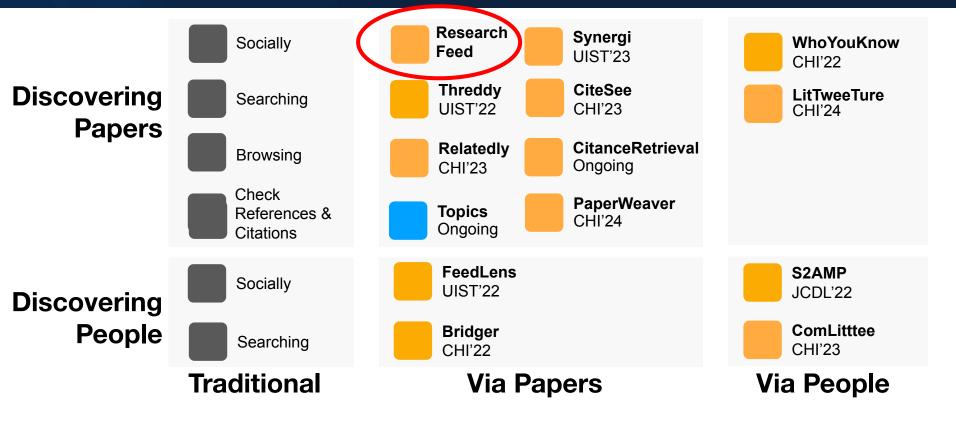


Discover!



Discovery - Our Research





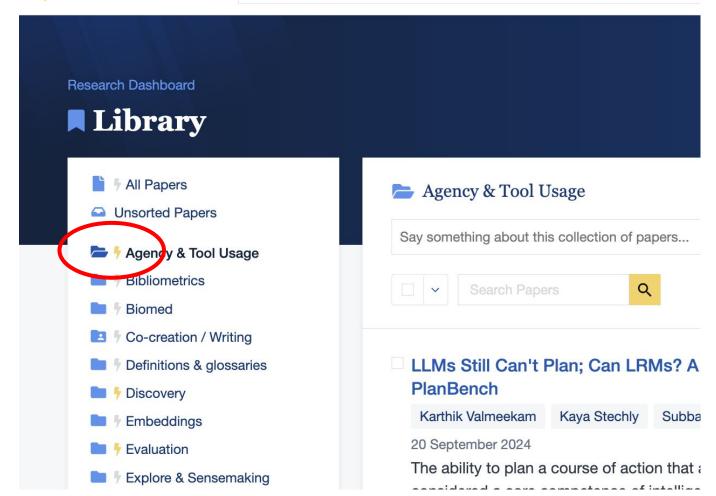


Adaptive Recommendations

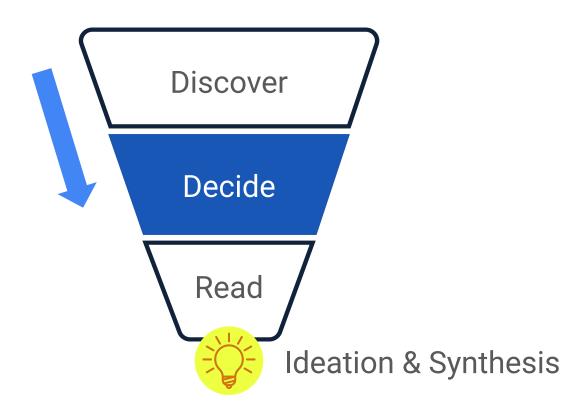


Papers you've rated, or browsed saved to library, or browsed S2 suggests new papers that user will find relevant 52 learns recommender **Previously Rated Papers** The Mythos of Model Interpretability (Lipton 2018) **Email Lung Disease Prediction** using ML (Monsi 2019)



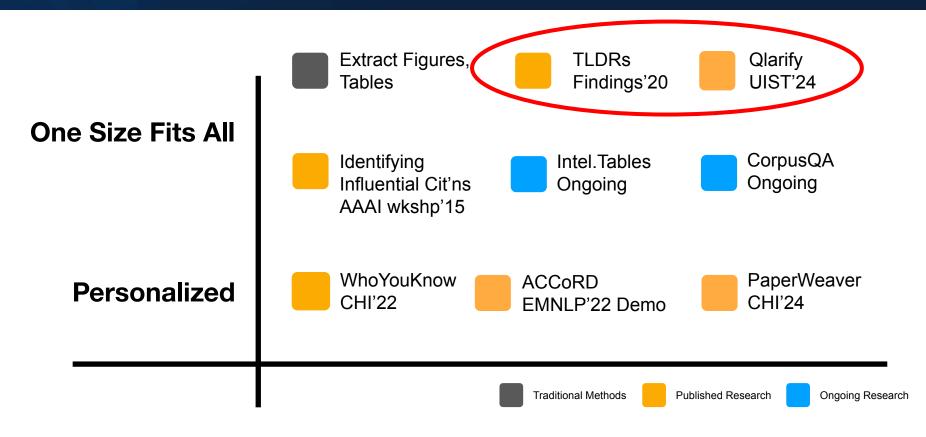


Decide!



Decide - Our Research





Extreme Summary TL;DRs









TL;DRs vs Abstracts

TLDRWe develop a new topological complexity measure for deep neural networks and demonstrate that it captures their

more tathomas

 \mathbf{Abs}

terrogating the network with input dan terizing and monitoring structural properties, however, have no been developed. In this work, we propose neural persistence, a complexity measure for neural network architectures based on topological data analysis on weighted stratified graphs. [...]

[...] In this work, we present the following contributions: We introduce neural persistence, a novel measure for characterizing the structural complexity of neural networks that can be efficiently computed. We prove its theoretical properties, such as upper and lower bounds, thereby arriving at a normalization for comparing neural networks of varying sizes. [...]

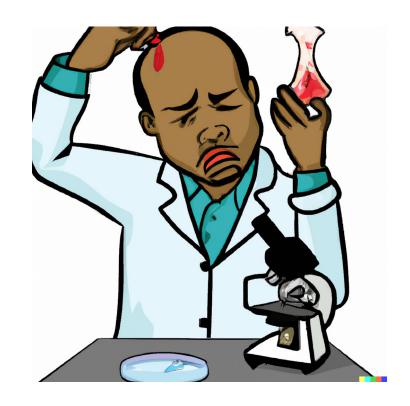
[...] However, this did not yield an early stopping measure because it was never triggered, thereby suggesting that neural persistence captures salient information that would otherwise be hidden among all the weights of a network [...]

How to Evaluate TLDRs



- Time on page?
- Click-thru rate?
- Bounce rate?
- User feedback:





Then... We added TLDRs to Email Alerts



- Tested: Impact of replacing truncated abstracts with model generated TLDRs in alert *emails*.
- Expected: An increase in clicks and saves.
- Hoped: Unsubscribe rate doesn't increase.

Dany Haddad. Medium, March 2023

https://medium.com/ai2-blog/tldrs-help-to-cut-through-the-clutter-3ad802caed93

Evaluating TLDRs in Alert Emails

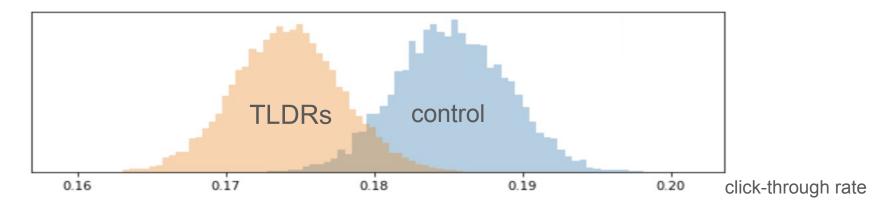


No increase in unsubscribe rate



• Stat significant 6% *decrease* in CTR



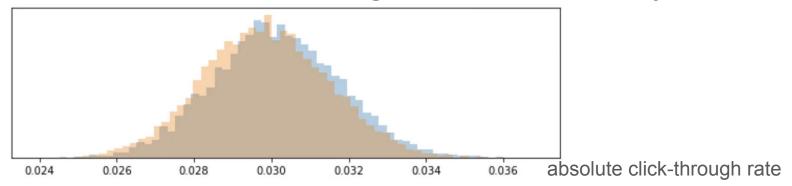


Why? Lower engagement? Greater efficiency?

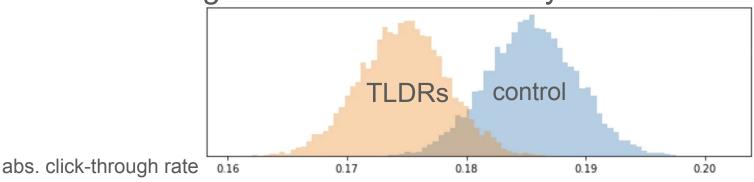
Click-thru Rate Conditioned on Action



Email → Click-through → Save to Library



Email → Click-through → **NO** Save to Library





Al-Powered Question Answering



Paper-Specific QA VS.

(asking about specific paper)

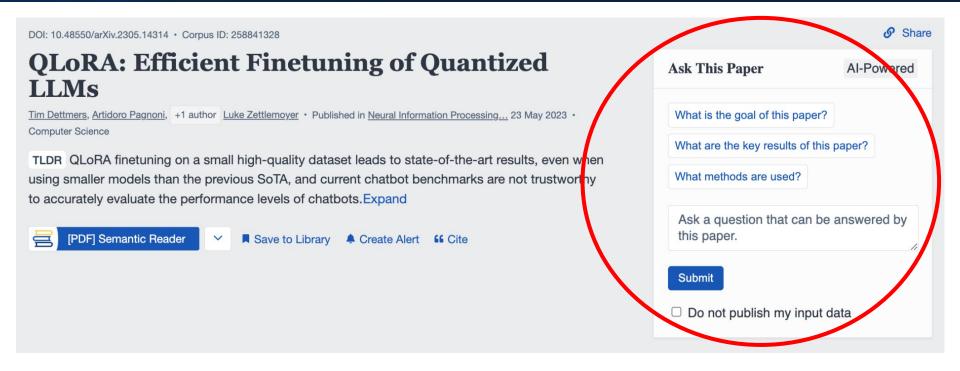


CorpusQA

(asking a question over millio of papers)



Paper-Specific Question Answering v1





Paper-Specific Question Answering







QLoRA introduces a number of innovations to s performance: (a) 4-bit NormalFloat (NF4), a nev

theoretically optimal for normally distributed weights (b) double quantization to reduce the average memory footprint by quantizing the quantization constants, and (c) paged optimizers to manage memory spikes.

OLoRA i perform theoreti reduce t and (c) r

CO

pa

Define "double quantization"

k-bit NormalFloat (NFk), since the data type is information-theoretically optimal for zero-centered normally distributed data. The exact values of this data type can be found in Appendix E.

Double Quantization We introduce Double Quantization (DQ), the process of quantizing the quantization constants for additional memory savings. While a small blocksize is required for precise 4-bit quantization [13], it also has a considerable memory overhead. For example, using 32-bit constants and a blocksize of 64 for W, quantization constants add 32/64 = 0.5 bits per parameter on average. Double Quantization helps reduce the memory footprint of quantization constants.

More specifically, Double Quantization treats quantization constants $c_2^{pq/2}$ of the first quantization as inputs to a second quantization. This second step yields the quantized quantization constants c_2^{pq} and the second level of quantization constants $c_1^{pq/2}$. We use 8-bit Floats with a blocksize of 256 for the second quantization as no performance degradation is observed for 8-bit quantization, in line with results from Dettmers and Zettlemoyer [13]. Since the c. 293 are positive, we subtract the mean from c2 before quantization to center the values around zero and make use of symmetric quantization. On average, for a blocksize of 64, this quantization reduces the memory footprint per parameter from 32/64 = 0.5 bits, to $8/64 + 32/(64 \cdot 256) = 0.127$ bits, a reduction of 0.373 bits per parameter.

Paged Optimizers use the NVIDIA unified memory 3 feature wich does automatic page-to-page transfers between the CPU and GPU for error-free GPU processing in the scenario where the GPU occasionally runs out-of-memory. The feature works like regular memory paging between CPU RAM and the disk. We use this feature to allocate paged memory for the optimizer states which are then automatically evicted to CPU RAM when the GPU runs out-of-memory and paged back into GPU memory when the memory is needed in the optimizer update step.

QLORA. Using the components described above, we define QLORA for a single linear layer in the quantized base model with a single LoRA adapter as follows:

$$\mathbf{Y}^{\text{BF}16} = \mathbf{X}^{\text{BF}16} \text{doubleDequant}(c_1^{\text{FP}32}, c_2^{\text{k-bit}}, \mathbf{W}^{\text{NF4}}) + \mathbf{X}^{\text{BF}16} \mathbf{L}_1^{\text{BF}16} \mathbf{L}_2^{\text{BF}16}, \tag{5}$$

where doubleDequant(-) is defined as:

$$doubleDequant(c_1^{FP32}, c_2^{k-bit}, \mathbf{W}^{k-bit}) = dequant(dequant(c_1^{FP32}, c_2^{k-bit}), \mathbf{W}^{4bit}) = \mathbf{W}^{BF16},$$
 (6)

We use NF4 for W and FP8 for c2. We use a blocksize of 64 for W for higher quantization precision and a blocksize of 256 for co to conserve memory.

For parameter updates only the gradient with respect to the error for the adapters weights $\frac{\partial E}{\partial E}$ are needed, and not for 4-bit weights $\frac{\partial E}{\partial W}$. However, the calculation of $\frac{\partial E}{\partial T}$ entails the calculation of $\frac{\partial X}{\partial W}$ which proceeds via equation (5) with dequantization from storage WNF4 to computation data type WBF16 to calculate the derivative $\frac{\partial X}{\partial X}$ in BFloat16 precision.

To summarize, QLORA has one storage data type (usually 4-bit NormalFloat) and a computation data type (16-bit BrainFloat). We dequantize the storage data type to the computation data type to perform the forward and backward pass, but we only compute weight gradients for the LoRA parameters which use 16-bit BrainFloat.

See in paper context

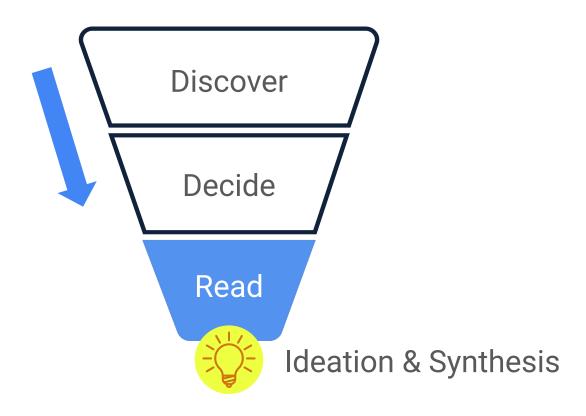
(DO), the process of ts for additional memory equired for precise 4-bit derable memory t constants and a constants add 32/64 = 0.5ble Quantization helps intization constants...

guantization onstants. The : memory

ribution!



Read!



Reading Technical Papers is Tiring!



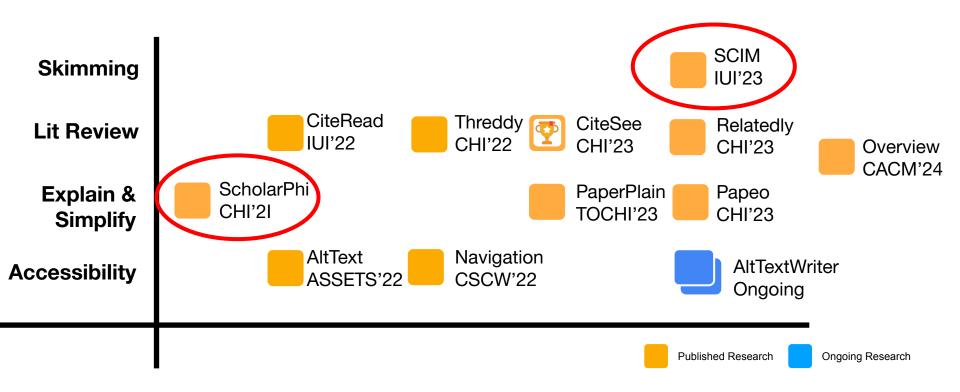
- Researchers flipped back and forth in document...
- Looking for definitions, checking citations...





Semantic Reader Project





Term & Symbol Definitions

Table 3. This difference in computation requirements is morequires an average of less than 5 seconds per alphabet whe 100. Table 4 additionally shows that in addition to being accurate in estimating k. This demonstrates the efficacy of cluster number (page 8). ■ ∑ ■ ×

Al-Powered Skimming



Goal

Enable skimming in less time, with less effort.

Challenge

How surface *salient* content w/ non-intrusive affordances?

Scim: Semantic Skimmer





Al-Powered Skimming



Figure 8. StateLens maintains a relatively stable error rate for state detection as the number of states increases, compared to the increasing trend in the baseline approach.

Number of states

increasing trend in the baseline approach (Figure 8). Next in user evaluation, we further demonstrate how the generated state diagrams power interactive applications to assist blind users access existing dynamic touchscreen devices.

USER EVALUATION

The goal of our year study was to avaluate how the components of StateLen agent, and people to a inaccessible

10 of 15

A) Highlight salient text

Annaratus and Participants

In order to enable repeated testing without wasting coffee, we built a simulated interactive prototype of the coffee machine in Figure 4 with InVision [23], which we displayed on an iPad tablet of similar size as the coffee machine's interface (iPad Pro 3rd generation, 11-inch, running iOS 12.2 without VoiceOver enabled). The conversational agent and the iOS polication were installed on an iPhone 6, running iOS 12.2 with VoiceOver enabled. The finger cap and the conductive stylus in Figure 3 were fabricated and used. We recruited 14 visually impaired users (9 female, 5 male, age 34-85). The demographics of our participants are shown in Table 3.

Procedure

Following a brief introduction of the study and demographic questions, participants first completed tasks using the 3Dprinted accessories. For each of the three screen placements (in the order of 90° vertical at chest-level, 45° tilted at chest-level. and 0° flat on the table), participants completed five trials using both the finger cap and the conductive stylus. The order of accessories was counterbalanced for all participants. For each trial, participants were first instructed to explore by placing the accessory on the touchscreen and move according to the

use the 3D-printed accessories to perform the tasks following the guidance and feedback of the iOS application. These realistic tasks involved a series of button pushes across many states, e.g., select gourmet drinks, cafe latte, strong strength, then confirm, auto-select default coffee bean, and end on the drink preparation screen. The iPad Pro simulating the inaccessible coffee machine was placed tilted at chest level, and the iPhone 6 running the iOS application was mounted on a head strap to simulate a head-mounted camera. Task completion rate and time were recorded.

After each step of the study, we collected Likert scale ratings and subjective feedback from the participants. Finally, we ended the study with a semi-structured interview asking for the participant's comments and suggestions on the StateLens system. The study took about two hours and participants were each compensated for \$50. The whole study was video and audio recorded for further analysis.

> ur user study results and summarize user ferences. For all Likert scale questions, parong a scale of 1 to 7, where 1 was extremely as extremely positive.

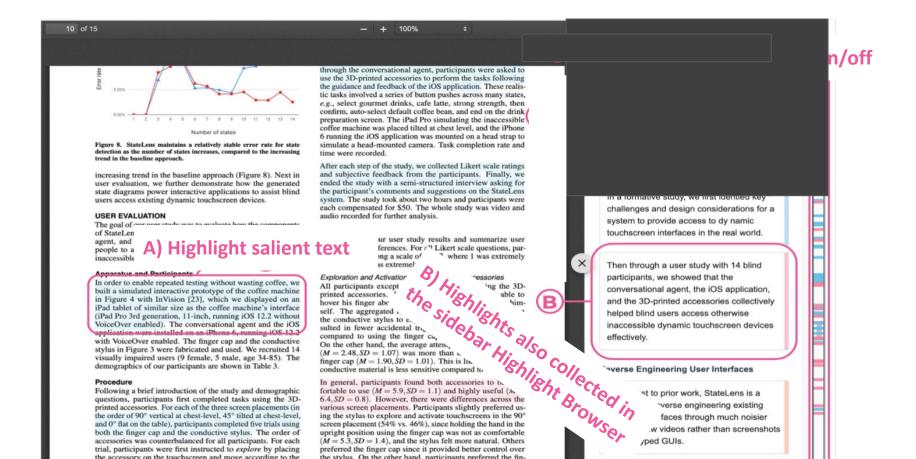
Exploration and Activation with 3D-Printed Accessories

All participants except P12 completed tasks using the 3Dprinted accessories. P12 had low vision, and was able to hover his finger above the target and then activate by himself. The aggregated results are shown in Table 4. Using the conductive stylus to explore touchscreens generally resulted in fewer accidental triggers (M = 0.03, SD = 0.16) compared to using the finger cap (M = 0.07, SD = 0.27). On the other hand, the average attempts of using the stylus (M = 2.48, SD = 1.07) was more than that from using the finger cap (M = 1.90, SD = 1.01). This is likely because the conductive material is less sensitive compared to fingers.

In general, participants found both accessories to be comfortable to use (M = 5.9, SD = 1.1) and highly useful (M =6.4, SD = 0.8). However, there were differences across the various screen placements. Participants slightly preferred using the stylus to explore and activate touchscreens in the 90° screen placement (54% vs. 46%), since holding the hand in the upright position using the finger cap was not as comfortable (M = 5.3, SD = 1.4), and the stylus felt more natural. Others preferred the finger cap since it provided better control over the stylus. On the other hand, participants preferred the fin-

Al-Powered Skimming





AI-Powered Skimming





Al-Powered Skimmingr





Al Helpful at Every Stage

Discover Decide Read **Synthesis**

Seek tools that include attribution & verification

Corpus QA techniques still developing

Thanks

Many co-authors and collaborators, but especially Raymond Fok, Marti Hearst, Andrew Head, Joseph Chang







