## **Optical Character** Software

actum diam est a felis. Na.. magna. dunt pu .s mus. F rcu interc la vitae, 1 euismod in elit tris aliquam a ndit eros, a n mi, eget llentesque to, in alique plandit et, l isellus acci indum a. ve a, id venena n. Sed mat. us. Nulla to taciti socios libero. Dor \* auque, po. Integer put toat. Duis ve.. ייי vida in vitae איים. Integer ייי "h pretium, liquia h

neque nec ci ad velit. Ir. er tristique c eget dui ut n risus eleit efficitur sce Curabitur n folor eget i et turpis qu liquet sap rra eget le natis susc on sagittis anim mat' au felis ia us varius bero inte udin ma' Jia nostr ∋ nibh e' uere vitr lisis alic Lintum sit

#### , porttite st. Orc Julvinar is velit t∈ ApplePickers August 8, 2018 la euismoc iquam effic

. Mauris qu in pulvinar ardum. Sed neque temp e, sed peller. Suspendisse hímenaeos Curabitur pell. toque penatibu. msan tortor, eu 📖 • quam. Curabitur ac felis dui. Nur. suscipit ut. Pellentesque \*\* amet sceleric:

Jonec r

ariu. a allus

rient

...a. Etian

h eu, ar

am volutp s quis nune umsan. Na te maona. tincidunt i r ridiculus ana quis arc isque et urna 🗸 anicula eros euismod po n ultriann

ue penatit

tortor, eu

am. Curab

is suscipit

irpis, sit ai

s. dui ut a

um. Fusce

sed dolor

ed fermen'

prem ipsur

us auctor (

m metus.

tiam pelle c turpis ar trius est r ) nisi non acus odir s varius malesur iam elementum, ultrices metur ar magna vitae tellus feugiat fa haretra in. E amet, cons is. Proin ur m est a fel. facilisis, vit hc sed diam. 'ui ornare v. sus, imperc modo puru ngilla tellus, Suspendiss. er rísus. Ut quis ta id lorem eu

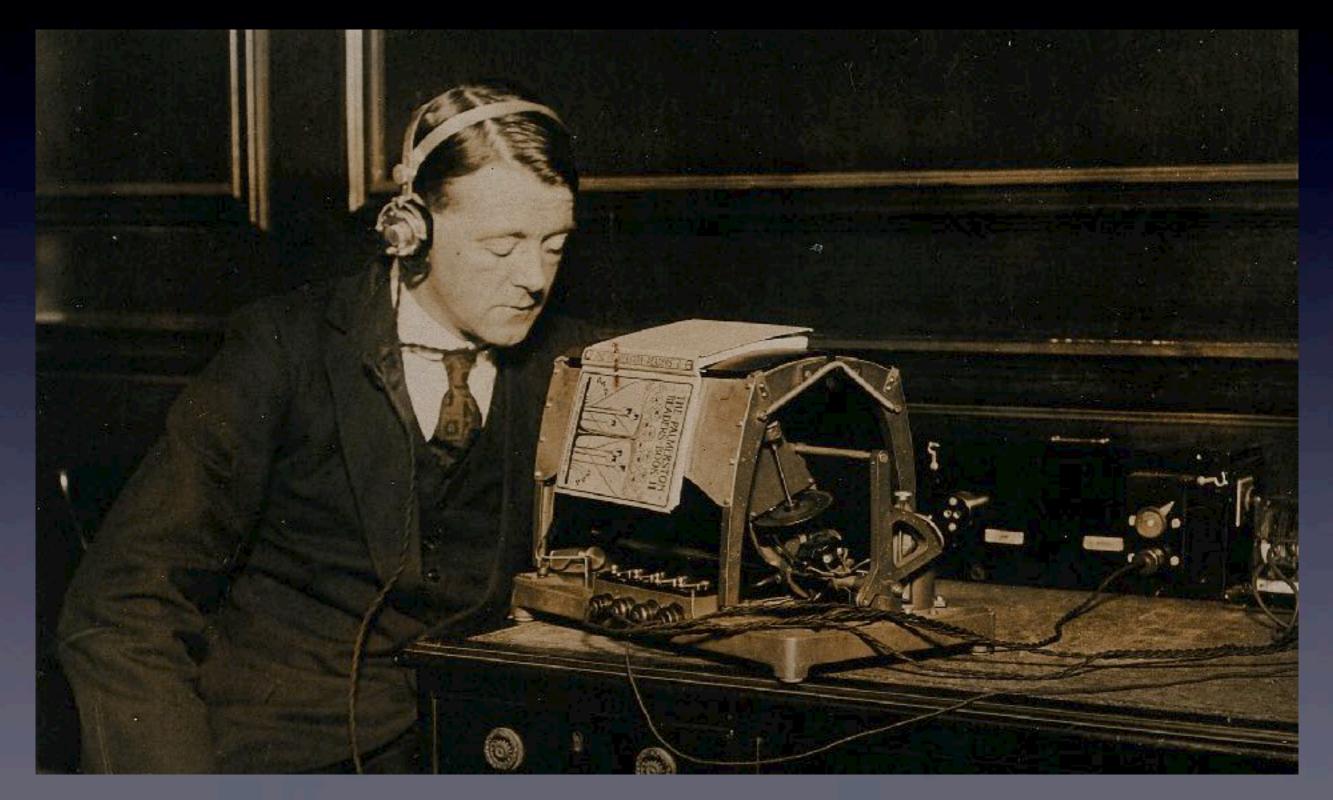
i volutpat. c

arturient m

# Brief History

- ✓ In 1914, Emanuel Goldberg developed a machine that read characters and converted them into standard telegraph code.
- Concurrently, Edmund Fournier d'Albe developed the <u>Optophone</u>, a handheld scanner that when moved across a printed page, produced tones that corresponded to specific letters or characters.
- ✓ In the late 1920s and 1930s Goldberg developed what he called a "Statistical Machine" for searching microfilm archives using an optical code recognition system.
- ✓ In 1974, Ray Kurzweil started the company Kurzweil Computer Products, Inc. and continued development of omni-font OCR, which could recognize text printed in virtually any font (Kurzweil is often credited with inventing omni-font OCR, but it was in use by companies, including CompuScan, in the late 1960s and 1970s
- Kurzweil decided that the best application of this technology would be to create a reading machine for the blind
  - Unveiled in 1976 and commercial product sold in 1978
  - Company sold to Xerox in 1980 which spun unit off as Scansoft which merged with Nuance Corp.
  - OCR to Braille available as OpenSource or <u>BrailleTranslator.org</u>

# Optophone



# Options

 Service companies offer bulk OCR to other companies

- Conversion of microfiche
- ExperVision

✓ Online free services
OCR Space

Name	Founded year	Latest stable version	Release year	License	Online	Windows	Mac OS X	Linux	BSD	Programming language	SDK?	Languages	Fonts	Output Formats	Notes
Tesseract	1985	3.05.02	2018	Apache	No	Yes	Yes	Yes	Yes	C++, C	Yes	100+[1]	Any printed font	Text, hOCR,	Created
Readiris	1986	16	?	Proprietary	?	Yes	Yes	?	?	?	Yes	100+[5]	?	?	Owned by Canon
Screenworm	2013	1	2014	Proprietary	No	No	Yes	No	No	Objective-C++	No	57	?	тхт	Product of Funchip. Uses
ExperVision <sup>[6]</sup> TypeRead er & RTK	1987	7.1.170.1125	2010	Proprietary	Yes	Yes	Yes	Yes	Yes	C/C++	Yes	21	2618		Has a Mobile and Embedded System version
ABBYY FineReader	1989	14	1/25/17	Proprietary	Yes	Yes	Yes	Yes	Yes	C/C++	Yes	192[8]	?	DOC, DOCX, XLS, XLSX, PPTX, RTF,	ABBYY also supplies SDKs for
Asprise OCRSDK	1998	15	2015	Proprietary	Yes	Yes	Yes	Yes	Yes	Java, C#,VB.NET, C/ C++/Delphi	Yes	20+[11]	?	Plain text, searchable PDF,	Java, C#, VB.NET, C/C++/Delphi
LEADTOOL S <sup>[16]</sup>	1990[17]	19	2014	Proprietary	Yes	Yes	Yes	Yes	No	C/C++, .NET, Objective-C, Java,	Yes	56[18]	Any printed font	PDF, PDF/A, DOC, DOCX, XLS, XPS,	Supports Latin, Asian, Arabic, and MICR
CuneiForm	1996	1.1	4/19/11	BSDvariant	No	Yes	Yes	Yes	Yes	C/C++	Yes	28	Any printed font	HTML, hOCR, native, RTF, TeX,	Enterprise-class system, can save text
OmniPage	1970s	19.2	2015	Proprietary	Yes	Yes	Yes	Yes	No	C/C++, C#[24]	Yes	125[25]	Machine and handprinted fonts	DOC/DOCX XLS/XLSX PPTX RTF PDF	Product of Nuance
gImageRead er[27]	2009	3.2.99	2017-07	GPL	No	Yes	Yes	Yes	No	C++	?	100+	Any printed font	TXT, PDF, hOCR	uses Tesseract OCR engine
GOCR	2000	0.5	2013	GPL	Yes[28]	Yes	Yes	Yes	Yes	С	?	20+	?		
Ocrad	?	0.25[29]	4/16/15	GPL	Yes	Yes	Yes	Yes	Yes	C++	Yes	Latin alphabet	?		Command line
SmartScore	1991	10.5.8	2015-07	Proprietary	No	Yes	Yes	No	No	?	?	?	?		For musical scores
PDF OCR X	2008	2.0.22	2016	Proprietary	No	Yes	Yes	No	No	Java, C++, Objective-C	No	100+	?	TXT, Searchable PDf	Drag and drop UI.
OCRopus	2007	1.3.3	12/16/17	Apache	No	No	Yes	Yes	Yes	Python	?	All languages using Latin	Normal Latin script	TXT, hOCR <sup>[30]</sup> , PDF <sup>[31]</sup>	Pluggable framework under
MathOCR	2014	0.0.3	2015	GPL	No	Yes	Yes	Yes	Yes	Java	?	?	?	HTML, LaTeX	Features mathematical formula
Yunmai OCR SDK	2002	1	2013	Proprietary	Yes	Yes	Yes	Yes	Yes	Java, C++, C, object pascal, objective-C	Yes	14	Any printed font	TXT, PDF	Has the advantage of
Anyline SDK	2013[34]	3.5.1[35]	2016[35]	Free non- commercial	No	No*	No*	No*	No*	Java (Android), Objective-C &	Yes[36]	2 (German, English)	Any printed trainable font[38]	Plain text, verification image	*Customizable mobile OCR SDK for
AliusDoc AD-SCI <sup>[7]</sup>	2005	2.1	2015	Proprietary	No	Yes	No	No	No	VB.Net	For Extensions	All ASCII- compatible languages	?	XML, PlainText, any other thru SDK extensions	Minimal need for post-sale Professional
E-aksharayan	2010					Yes	No	Yes	No			14		RTF, TXT, BRL	
Nicomsoft OCR SDK	1999	5.5	2015	Proprietary	No	Yes	No	Yes	No	C#, VB.NET, C+ +, Delphi, Java	Yes	25+[14]	?	Searchable PDF, Text, RTF	C#, VB.NET, C++, Delphi, Java OCR
AnyDoc Software	1989	?	?	Proprietary	No	Yes	No	No	No	VBScript	?	?	?		Works with structured, semi-structured,
OCR.space	2015	3.02	2017	GPL	Yes	Yes	No	No	No	C#	Yes	23	Any printed font	тхт	Windows desktop software,
SimpleOCR	2002	3.5	2008	Proprietary	No	Yes	No	No	No	?	?	?	?		
Dynamsoft OCR SDK	2003	8.2	2012	Proprietary	Yes	Yes	No	No	No	C/C++	Yes	40+[23]	?	PDF, TXT	
Microsoft Office	2011	?	2007	Proprietary	No	Yes	No	No	No	?	?	?	?		
FreeOCR	?	4.2	Aug-12	Proprietary	No	Yes	No	No	No	?	?	?	?		[26]
Microsoft Office	?	Office 2007	2007	Proprietary	No	Yes	No	No	No	?	?	?	?		Uses OmniPage[citatio
OCR.net	2016	?	2016	Proprietary	Yes	No	No	No	No	Java, C++, PHP, Objective-c	No	100+	?	TXT, Searchable PDF	Online service powered by PDF OCR X for
Puma.NET	?	?	10/29/09	BSD	No	Yes	No	No	No	C#	Yes	28	Any printed font		.NET OCR SDK based on Cognitive
ReadSoft	?	?	?	Proprietary	No	Yes	No	No	No	?	?	?	?		Scan, capture and classify business
Scantron	?	?	?	Proprietary	No	Yes	No	No	No	?	?	?	?		For working with localized interfaces,
OCRFeeder	2009-03	0.8.1	12/22/14	GPL	No	No	No	Yes	No	Python	?	?	?		Features a full user interface and has a
MeOCR	2012	1.0.0	2012	Freeware	No	Yes	No	No	No	C/C++/C#	Yes	28	Any printed font	HTML, hOCR, native, RTF, TeX,	Windows application.

### Pros and Cons of Macbased applications

- $\checkmark$  Some have been around for a very long time
- $\checkmark$  Most are automatic with little capability to tweak output
- $\checkmark$  Generally, you get what you pay for
  - None do a mistake-free job
  - Output should be reviewed
- Image resolution, font size, line straightness can dramatically affect quality
- $\checkmark$  All can handle multiple languages
- ✓ Most all use either OpenSource or proprietary 3rd party OCR engines such as Tessaract which is now Google's ML OCR Kit

Search Results for "OC8" 1-120 of 172 Sert By: Relevance PDEScanner - Scannin HP Easy Sean PDF Office: Edit Test PDF OCR X Communit. Presenter PDF  $\odot$ Productivity Utilities Rusiness. Productivity. Utifies 常常常常客 185 Belings 青吉吉宮宮 Sai Patings 吉吉吉宮宮 42 Balings OPEN -514.69 = 50.99 -DET -OPEN -IText - CCR & Translat. PDF Converter Pro OCR App by LEADTO. PDE-OCR-Eree PDE-to-Word-Free Productivity: Productivity Utilities Productivity Business 3 ★★★★☆ 15 Redings ★★文文文 24 Eatings 青青青岩岩 14 640 654 ★★☆☆☆ i6 katings 古古古古古 is hadingal \$38.92 -GET . CET -DET -021 -In-Appa Purchases. In-App Purchasea. In-Aco Purchases. PDF to Word with OCR. PDF OCH X Enterprise... PDF Converter with O., Mathpik Snipping Tool FineReader OCR Pro PDF 🤐 Productivity Productivity Business Productivity Productivity. ★★★☆☆ 37 Ratings ★★★おお 40 Radings オキオネネ 6 Radings ★★★☆☆ 5 Radings 812.65 -524.06 -DET -6110.05 -523.50 -PUF OCRKIL Prizmo 3 - Scanning &. PDF-element 8 Pro - P., Condense Tessa - OCR Productivity Business. Productivity. Utilities Preductivity ★★★☆☆ 5 Ballings 🔺 📩 📩 🖄 20 Ratings 青素素素論 B Ratings ★★★白白 A Ballings Abc OCR \$25.00 \* \$36.90 T 86.00 = SFT T 6-1 7 In-App Purchases In-App Purchases In-App Purchases PDF-Converter-Free OCRTOOLS Sign Master - PDF Ass... PDF-to-Excel+Free PDF to Excel OCR Con., Productivily. Productivity Precluctivity Numiroses, Rusin-sc **ABC** PDP 🚖 🖄 🖄 🖄 🛽 Ratings ★★☆☆☆ 17 Ratings ★★素☆☆ 10 Ratings ★★治疗治 s Ratings BET T BET T 661 7 BET T \$4.58 1 In-App Furchases In-App Purchases In-App Purchases PDF Converter OCR PDF Text Extractor - F. DecuSean Plus PDF PDF Editor Pro-PDE-to-Pycel-Pro Business. Productivity. Buy mean Productivity. Business DsP ★★★★★ 10 Patheor ★★★☆☆ 31 Willings ★★☆☆☆ 9 Redinps ★安安安安 5 Balings Tert 601 -953.95 -65T -34.59 -902.00 -In-App Purchases Text Extractor - Extra-PDF to Fucel with OGR PDF-Scanner-Pro Paperless. PDF to Word ±t PDF TEXT Business Business Business: Productivity: Businese ★★★★☆☆ 11 Racings ★★☆☆☆ Sitedings ★★★★☆ BO Ratings オ市市市市 7 Radings 512.29 -0 ST = 1 52.29 -538.02 = 515.98 -In-Acc Parchases In-Aco Purchases PDF-to-ePub-Free PD--to-HTML-Free Aristocrat PDF & Image Text Extr. Visualizer Productivity: Utilities Productivity Utilities Business OCR ★☆沖☆沖 e katings ★★☆☆☆ Billiadings ★★★★☆ S Bailings ★★★★☆ | Fatings TEXT DPT -6.97 -OFT -84.00 -0.07 In-App Purchases In-App Purchases PocketDrone [Service] Image-to-PDF-Free PDF to Text with OCR PDF-to-PowerPoint-Fr. DocScanner Productivity. Business Preductivity. Preductivity: Preductivity: ★★☆☆☆ 18 Ratings 🔺 🕸 🚖 🚖 🛚 Rallinge 99 6FT # 115.59 1 SFT T 8790 -In-App Purchases In-App Purchases eDec Organizer Cloud. PDF-to-Text-Free Shimo Scanner - PDF ... ExactScan Pro + PDF Converter OCR PDF Productivity. Productivity. Productivity Rusiness Maniputer. 0 青青青青宮 Billalings 919.59 921 -993.09 -95T -912.59 -In-App Purchases PDF Fdit Express - PD. PDF Foitor Plus - PDF PDF Toolbox + Easy Screen OCR Scanner Pro -PDF Doc. Productivity Productivity B. piness Productivity. Productivity ★★★☆☆ 5 Bedrigs ★★★★☆ So liantings PDF 54.53 = 36.39 30.99 -519.99 -FICOneers3 BET 📼 OCRWizard - Convert . WorldCard Subtitle Extractor PDF to HTML with OCR PDF Converter ++ POF POR Utilities utilities Business Video Productivity OCR ★★★☆☆ 10 Radings 古古古古古 17 Bacings \$53.00 -6.ET -54.60 -HTML 57.30 -525.05 -

1

# Two Major Groups

Those that come with scanning software or scanner itself, and those that do not

• Example: HP Easy Scan

Get from Mac App Store



✓ Generally those that focus solely on OCR do a better job

#### Stand-alone Applications/ Predominantly OCR

	<u>Fine</u> <u>Reader Pro</u>	<u>Prizmo</u>	<u>Cisdem</u> OCR Wizard	<u>PDF OCR</u>
		Abc Abc		
MSRP	\$120	\$50	\$60	\$0
Pros	Best accuracy, handles tabular data, scans	Multi-column text, Business cards	Easy to use	Free, simple interface
Cons	Cost, erratic support	Can't handle tabular data	Poor with tabular data	Output only PDF or .txt files

#### **Combination Products**

	<u>Acrobat</u>	<u>PDFPen</u>	<u>VueScan</u>	
		San La		
Cost	\$15/mo	\$75/\$125	\$40	
Pros	Best quality, many other features	Mac Centric, other features	Works with almost any scanner	
Cons	Beyond expensive	Good OCR, but exports file	Average OCR, complicated layout	

# OCR Everywhere

- There are several websites that offer OCR services
  - Minimal to no tweaking of scan
- ✓ OCR is built in to several iOS apps including Adobe Acrobat, MS-OneNote, Text Scanner and Office Lens