

5. 5, 구-본-
t로a 지u- 전a지
대 대 대 대
수민 v등배
람.

람-은ㄹ · ㅁ
지 주uuㄹ · ㄹ
tk지 동ak지 .
5.

tesseract chosen-image-1962-1973.jpg out -l kor
Outputted text from the combined trainneddata ran on
the picture on page 2

달, 달, 무슨 달.
낫 파 같이 밟은 달.
어디 어디 비추나.
우리 동네 비추지.

달, 달, 무슨 달.
거울 같은 보름달.
무엇 무엇 비추나.
우리 얼굴 비추지.



Potential titles.

타이포잔치 2013
서울 국제 타이포그래피 비엔날레
OSP 워크숍
9월 30일 - 10월 4일
문화역서울 284

An OSP workshop
at Typojanchi 2013
30/9 - 4/10
Seoul Station

We've hesitated a lot, finally not to decide a title for the workshop.

Optically Recognizing Characters?
Automating observations?
Human reading machine?

- Declutching OCR
- Can we read like machines?
- Protocr

Or from sentences of Tesseract
training documentation
— Character normalization
sensitivity prototypes
— Could this output normproto?

Or 'simply'
— Normprotocr
— Normproto

Or elsewhere
— Each letter is replaced by a ?
— How little can you get away with?
— Declutching training data
— How little can training data be?
— OCR training data
— Training data for OCR
— Can training data be reclaimed?

Description of the workshop as announced

Fancy reading machines and androids from science-fiction fantasies are embodied in our modern lower-profile world as OCR software packages. OCR means Optical Character Recognition and it is software that can extract text from image files. One of these softwares, the free and open source Tesseract¹ is composed of two parts that we can study, thanks to its license. There is the engine itself, and the training data for a language² partly based on what Tesseract called 'prototypes'. We could compare this 'before the type' (proto-type) to the culture a lecturer progressively gathers from his first lesson going from a novice to a fully grown expert. By following the limits between the blank surfaces and the dark pixels of the shapes of letters, Tesseract compares its journey with previous ones, on images already followed in the past. It starts by learning patterns and specificities of languages, rhythms and irregularities. It goes on to recognise the body of a glyph, then it works out, bit by bit, if this glyph is a letter, form is a word, and eventually it makes out phrases.

1 - [http://en.wikipedia.org/wiki/Tesseract_\(software\)](http://en.wikipedia.org/wiki/Tesseract_(software))
2 - <https://code.google.com/p/tesseract-ocr/>

Like all of us, Tesseract learns typography in this same process, in a completely intertwined way, as sentences, script and eventually, language.³

Tesseract follows rules by which it can make decisions. In a basic example from Latin script, if the software seems to be recognising something resembling to iii (three times the letter 'I'), specific rules kick in to suggest that it is most lightly the letter 'm' and not a triple consonant. Grammar and language coming in at a later stage, as it did for us, still following this unusual idea of teaching software to read.⁴ The very specificities of typography and how each shape is drawn and could or couldn't be distinguished from another one arrives just after. As in the previous example the potential small parts that protrude from the I are more likely to be the arc of the m if the font is a serif one than it is a sans serif one. This process becomes intertwined with the actual context: with time, the system becomes familiar, and extremely efficient with some specificities of a typeface. Its shape, its overall form and size now mean something. It would have to relearn an entirely new toolkit to be able to read a different typeface. With this, could the relations binding shapes to their meanings be noticed?

At young, naive and early stages of deciphering writing systems, slowly working out the building blocks to a

3 - <http://code.google.com/p/tesseract-ocr/wiki/TrainingTesseract3>

4 - [http://code.google.com/p/tesseract-ocr/wiki/TrainingTesseract3#The_last_file_\(unicarambigs\)](http://code.google.com/p/tesseract-ocr/wiki/TrainingTesseract3#The_last_file_(unicarambigs))

legible language, we wonder how synthetic constructions (like Hangul) compare to agglutinated ones (like Latin). More specifically, how do these methods influence OCR data? On a more contemporary note, it would be hard to deny how much screens and screen text technologies have influenced typography these days. All languages carry different meanings, different cultures with their characters. These grid(ty) displays are no favour to typographic heritage, but they have brought on so much interesting conundrums. The rendering engine ttf autohint voluntarily distorts vector shapes of glyphs to optimise screen rendering⁵. In this workshop, we propose to carefully replay some of the processes the OCR system uses to reread typography from the departure point of any new learner, the one we all have known at first and mostly definitively forgotten by now... By patiently observing the various parameters at play when a letter is to be differentiated from another, the thin and variable line of separation between signification and shape, between letter and typography begins to reveal itself. Could the different parts of the letters that compose barebones of other letters be recreated in a kind of wild reverse engineered Metafont⁶ paradigm, where all of the shapes of the glyph are defined with geometrical equations?

We wonder how much we can learn from methods borrowed off OCR. By replaying its methods, but basing

ourselves on some parameters only, not aiming for full comprehension, but basic knowledge of how our different sets of characters work retracing its first steps only? Would the outcome of this be enough to go on to understanding typographic subtleties, enabling a bridge between specificities in shape and specificities in language? Finally, if we know organisation in Hangul and Latin are different, and that they do work along with similar ideas, could we try to avoid the main caveats of forcing comparisons between each? Instead can we focus on the systems that the OCR-by-human must use to read both for rethinking deeper specificities between the two composition methods, between these two typographies, between these languages?

Description from typojanchi.org

OCR(광학 문자 인식) 장치의 작동 원리를 연구하고, 문자로 정제되기 이전의 타이포그래피 형태를 관찰해보는 워크숍. 끈질긴 리버스 엔지니어링을 통해 타이포그래피를 더욱 깊이 이해하고, 그 형태와 의미의 관계를 되새겨본다.

강사:
OSP (사라 마냥, 피에르 하위허베르트,
뤼디빈 루아소, 콜름 오닐)

강사 소개

5 - [http://www.freetype.org/ttfautohint/
#samples](http://www.freetype.org/ttfautohint/#samples)

6 - <http://en.wikipedia.org/wiki/Metafont>

Workshop process

Monday and Tuesday were dedicated to OSP and workshop presentation. We based our talk about OSP on project and workshop that were linked to the Typojanchi workshop.

Virtual machines were set up and the participants had their first steps into this environment, combining the command line, paths to files, and exploring new tools. We finished off the day with a dry run of the Tesseract program, our first way into Optical Character Recognition.

We built image samples on Wednesday that brought us to the first steps of the Tesseract training process.

From there, the workshop process and the chronological steps we had planned, went out the window and forced us into a nice switch.

We needed to find a way to customise our virtual machines to enable Korean character writing. This was essential for the first steps of the training process during which we match coordinates on a picture to a unicode character.

The menus and all of the interfaces on the machines were switched to a language that none of the OSPs could understand, obliging the participants to explain what was happening in the menus, dialog, etc.! **The language and dual alphabets issue was definitely a good way of forcing the roles to change, flipping skills with knowledge!**

59 불은 수본

불은 자본 권불수할

식 민 민 민

수광 등민

할0

음수은

번 주번등

번 주자 0

할화

№ 006693 왕
승 선 권

청평사
오봉산 관광지 ← 댐

요금 ₩3,000원

소양관광개발주식회사

본권은 환불할 수 없음

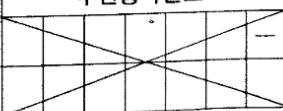
№ 006693 복
승 선 권

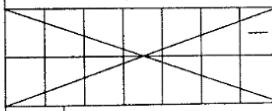
청평사
오봉산 관광지 → 댐

요금 ₩3,000원

소양관광개발주식회사

본권은 환불할 수 없음

승 선 자	
성명 :	
주민등록번호	
	
전화 번호	

승 선 자	
성명 :	
주민등록번호	
	
전화 번호	

006693 요

은은은 006693 원

승 선 권

산수권
호등청 55지 ←댐

요금 ₩3,000원

광소양권광개발주식회사W

사등권본 환사등환 수 없등

소사민관
불등민발

청평사 관
오봉산 복광광등광변댐

요금 ₩은,000원
소양권광광권발주식회사

본권은 환불할 수 없등

۔۔۔ ?স, ম.
ৰ ৰ প প ন ৰ ৰ ও প
০ । ০. ০ । প ।
-ন-মপ -ম ম ম ম

. ৰ .

?ৰ প ল - ? প
. ম ৰ স প ম ম প
প ম ? প স ম .

. ৰ .

“나는 그에게서 모든 것을 배웠다.”
“나는 그에게서 모든 것을 배웠다.”

4-1-4-3

UU UU UUUU
U UUU UUU
UU UU UU UU
UUUU UUUU

UU

UUU
UU UUUU
UU UUUU U
UU



u보대 tuu보

도보FF보대내M과보 괴보F대보v

M도MFM 도

과보

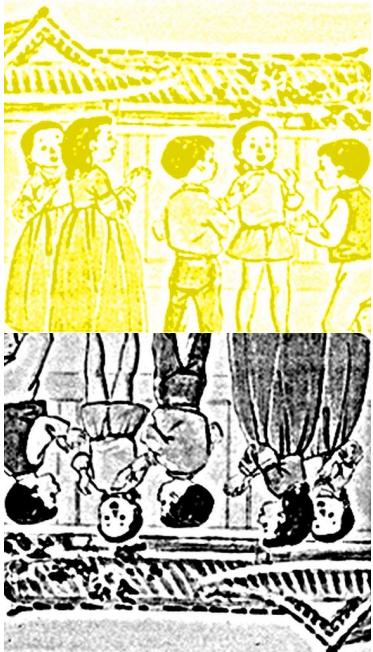
보보 보과괴면q괴과q보k q|보보괴과
대 보kk괴 보괴보
보보면괴 대되보 보도대 k보되보과

보대내과보 되Mk보 보t보q보q 내보
보

괴q보보 되보 k보면보내 괴보 보도괴
괴보v보대보a 대보MFM괴 보괴보g
보되a괴q보보보보. 괴보내va내보
괴보보보보q보 F면utvt면u면u
내 내MF과F과되대보 과Fkv보보qv
보보보과F보보 보FM보보과보 a

괴과a도보보보보보보 보대내
aFaFaFavav a보g대 t보과보과대 보
보보 보보보 괴과v보보내과a보F과보k
내괴v대vk보 괴보보Ft보보보도보
kvvt도 보보 보보과 보보보 보보kF
보 보 도과보 보보보g F보과보k과
내대 보보 보보k과 괴과M 보괴k보 q
보보k되보과

흰 달 흰 달 무슨 달.
 짜 오 짜 오라 흰 달은 데이
 거울 같은 달에 린이.
 어 어 어 라
 등구엇으다 털리 바빠나.
 달은 리 얼굴 비추지.

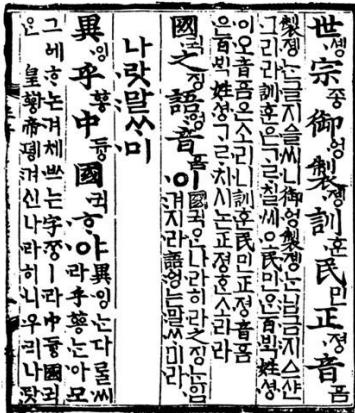


고는 놀 놀 놀고나.
 고는 놀 놀 놀고나.
 고는 놀 놀 놀고나.
 고는 놀 놀 놀고나.

달, 달, 무슨 달.
 낫과 같이 밝은 달.
 어디 어디 비추나.
 우리 동네 비추지.



통통한 통통한 통통한 통
 통통한 통통한 통통한 통
 통통한 통통한 통통한 통
 통통한 통통한 통통한 통



세종대왕 Sejong the Great

컴퓨터를 읽다
computer reads



나랏말쓰·미 등국에달아
문쓰·○와로 서로 스·ㅁ·스디 아니흐·ㄹ쓰··]
이런 편쓰·로 어린 ㅂ·] 그 성이
니르고저 호从此이서도
ㅁ·ㅊ·ㅁ내 제 ㅍ—들 시러퍼디
몬흐·ㅋ 노미 아니라
내이르·ㄹ 임흐·야 어엿비너겨
새로 스를여들쓰·○ㄹ·ㄹ·ㅁ·ㅓ(앵)ㄱ·노니
사르·ㅁ마다 ㅎ·ㅓㅓㅋ 수녕 | 니겨 날로새—메
뻔히 허느끼 ㅎ·고져 ㅎ·ㅋ 쟤·(마) ㄹ·미니라

달, 달, 무슨 달.
낫 파 같이 밝은 달.
어디 어디 비추나.
우리 동네 비추지.



kor

.탠 내폈
야!] 흠 주주
오 베 미
슨
우 은 앗 굴
갈 우얼
딴
.움 엇리
달 거 우 우
탱
.온냐:샤
닫 맑 추 추
바 바
슨 이
우갈커뻬
꾀뻬챤
닫꽈°
, 이 리
닫 낫 어 우

달, 달, 무슨 달.
거울 같은 보름달.
무엇 무엇 비추나.
우리 얼굴 비추지.



eng

'cf. '2, -'?-E-
'5!' 3l-- '2'1
019] of P]
—?-al -241]

'ii.

'H-8:
"I \$1'!-.
"I 337' .

'3.

일민미술관
ilmin museum of art

서울시 종로구 세종로 139번지 110-050
139 Sejong-no, Jongno-gu, Seoul 110-050, Korea

Telephone 82 2 2020 2050 | Facsimile 82 2 2020 2069
www.ilmin.org



정시관람권
Exhibition Ticket

No. 22207

kor

01 꾀 짹량 똑
꼰 같 미 놀 난
줌|m줌|1 m|[8순 [|m 0{ 컨빼
서울시 종로구 세종로 139번지 110-050
139 86] 0꼬담-꼬0,] 0꼬등꼬0-담맴,
86001 110-050, 1 <0> '6웹

[`(珥16[]) 0[] 6 82 2 2020 2050 1
줌웹:5珥m116 82 2 2020 2069

밖빼~크1m코[].0 [등

페흐

전시관람권
묘×뇨썰빼0꼬 쌈뿡뇨값

페 22207

eng

O] I] * 'V
a 1.] Bl E '3
ilmin museum of art
/4%/-~] %\$:;7- /~']]%\$ 139'.51_Z]
110-050
139 Sejong-no, Jongno-gu, Seoul
110-050, Korea

Telephone 82 2 2020 2050 I
Facsimile 82 2 2020 2069

www.ilmin.org

.. /

\$1/~15!-} 'li'4_
Exhibition Ticket

22207

달, 달, 무슨 달.
낮 파 같이 밝은 달.
어디 어디 비추나.
우리 동네 비추지.



달, 달, 무슨 달.
겨울 같은 보름달.
무엇 무엇 비추나.
우리 얼굴 비추지.



상상 의실상 의의실의 상상상의실의 실의
의의 의의의 상상 의상의 의의의의
상 상 상상상 상 의상상
상실의의 실상상상



정의상실

정의상상실

명동의상실

상상의상상

노영의상실

정의상상의상상

서울의상실

상상의상실

Genealogy

A genealogy of several OSP projects unfolds since 2008 around the questions we'll be excited to explore in another way in Seoul :

— During several Dingbat Liberation Fests, OSP insisted on questioning from the side what a glyph is, what a letter is, and how we can reappropriate these precious parts of the culture as our own (even if a legitimate smoothing effect by Unicode codification effort tend to loose the granularity at the heart of each script-language-typography triangle). Dingbat Dictées are recipes for teaching Unicode poetry to students of all ages.
<http://ospublish.constantvzw.org/blog/?s=dingbat>
<http://www.constantvzw.org/vj12/>
<http://ospublish.constantvzw.org/blog/typo/dingbat-dictee>
<http://ospublish.constantvzw.org/blog/live/dingbat-liberation-fest-ii>

— An osp installation system called 'Nancy', deploying the Dingbats Liberation Fest at the gallery My.Monkey in Nancy, (F) and then at Make-Art Festival in Poitiers, (F).
<http://ospublish.constantvzw.org/blog/news/dingbats-in-a-monkey>
<http://ospublish.constantvzw.org/nancy/>

— Based on these experiments, OSP member Pierre Marchand develops 'Nancy', a software made to support collective font design, also if dispersed over time and geographic location. A patchwork of classic Free, Libre and Open Source softwares, the system is

semi-automated process of scanning of cardboard-cut dingbats and a collaborative font.

— Part of Nancy is the OCR module 'Fonzie': Fonzie outgrowth has been developed after several pushes between, among others things, the need for a convenient way to reproduce handdrawn fonts for translated version of comics, using Opentype features to automatically switch between different versions of each glyph that try to emulate variations of typical hand drawn lettering.
Fonzie's progressive modifications follow various different type design needs and shapes multiple processes.
<http://ospublish.constantvzw.org/blog/?s=fonzie>
<http://ospublish.constantvzw.org/blog/news/asian-record-side-a>

— A more straightforward version of Fonzie is used in workshops starting schools in Brussels and London, but also performances and book design.

— A hackerspace in Brussels built a DIY-bookscanner and another evolution of Fonzie nicknamed Funzie added other OCR features to convert physical to digital reading by producing hybrid fonts. The reading process of the machine can be watched, understood, read, rewritten, changed and endlessly executed. This explores the artistic potential of the different elements in the transmission from the physical book to the digital object: the code, the fonts, the training data, the book and the digitally reborn text.
<http://video.constantvzw.org/VJ13/>
http://www.vj13.constantvzw.org/arc_hive/funzie_fonzie/

Exchange

Ray Smith is the main developer of Tesseract since 1985.

Ray Smith
<theraysmith@gmail.com>
1 Oct 2013
to Colm, Sarah, Pierre, Ludi

The training process that we use involves rendering text in available fonts and building a language model.

The source text for both these processes is gathered from web pages that has been identified to be in the required language.

Some languages have received special treatment for specific problems, but Korean is not one of them.

Although people occasionally ask for the training data, it is just too big to host on the site even if we had all the necessary copyright clearances.

Instead we are working on opening up the training tools that we use so more people can enjoy automated training.

Hope that helps.
Regards
Ray.

Sep 28, 2013, Colm O'Neill:

Dear Ray Smith,

We're contacting you regarding Tesseract, some of the dev forums suggest you as a person to contact about language support in Tesseract.

We are OSP (Open Source Publishings), a working group based in Brussels, that works exclusively with Free / Libre Open Source Software in the context of Graphic Design, bringing together commissions, teaching and research.
<http://osp.constantvzw.org/>

We're currently in Seoul prepping for a workshop at the typography biennial Typojanchi and for this we're going to use Tesseract as common ground between our different alphabets, i.e. Latin and Hangul. Our focus is going to be on the training of Tesseract, going through the steps to build some new specific data for the Korean language.

With this, we're wondering about the origins of the current data available for Korean. We can't find out where the downloadable data originates from, on what kind of text images it was trained, etc.

The history of the tool obviously has something to do with this, so we're coming to you hoping you can provide us with information about the origins? (What are the /) Where could we find the original samples?

Any help would be appreciated.

(You're not in Seoul yourself next week by the wildest of chances?)

All the best,
Yours,

Sarah, Pierre, Ludi, Colm,
for OSP.

Tools

This workshop was made possible
with the following free tools:

Tesseract-ocr 3.02

Git

Gimp

Inkscape

Libre Office

To help us generate good Tesseract boxes, we've used
<http://pp19dd.com/tesseract-ocr-chopper/?i=ocrJd0HbJ>

We have used also
an older Tesseract version -
svn checkout -r 659
<http://tesseract-ocr.googlecode.com/svn/trunk/> tesseract

Most of us were using Ubuntu
running on a Virtualbox.

The booklet has been designed in
Libre Office and Inkscape by the
participants, using Dejavu Serif
font. <http://dejavu-fonts.org>

The pdf pages will be compiled
through Podofo Impose if it works -
or we will impose it by hands with
Scribus!

A big thank you to the authors of all
these tools !

Biography

Open Source Publishing makes graphic design using only free and open source software—pieces of software that invite their users to take part in their elaboration.

Founded in 2006 in the context of Brussels art organisation Constant, the OSP caravan now comprises a group of individuals from different background and practices: typography, graphic design, cartography, programming, mathematics, writing, performance. Through a collaborative practice, they work on workshops, commissioned or self-commissioned projects, searching to redefine their playground, digging towards a more intimate relation with the tools.

True to their name, OSP publishes all the the source files to their projects through their website <http://osp.constantvzw.org>.

Sarah Magnan

Graphic designer, Sarah started to experiment in ERG (Brussels) possible links between graphic design and new media art. From links to links she became curious and interested by collaborative work, sharing matters on web, on print and more widely on archiving matters: which status to gives to archive, how to make it born or reborn, how to share it, show it, confront it.

Pierre Huyghebaert

Exploring several practices around graphic design, he currently drives the studio Speculoos. Pierre is interested in using free software to re-learn to work in other ways and collaboratively on cartography, type design, web interface, schematic illustration, book design and teaching these practices. Along with participating in OSP, he articulate residential spaces and narratives through the artists temporary alliance Potential Estate and develop collaborative and subjective mapping with Towards and others Brussels urban projects.

Ludivine Loiseau

Formed at the Ecole Estienne (Paris), Ludi relearned everything in Brussels. She immersed herself in the centre for graphical delicacies Speculoos and met the OSP group aboard a van on route to Poland in 2008. Ludi questions the contemporary role of typography and her practice is reflected in her courses at the Ecole de Recherche Graphic, where she is a lecturer in typography and free software. Ludi also works with Mathieu Gabiot on advancing the issue of licensing furniture objects, and supports an ephemeral publishing project, Le Calendrier.

Colm O'Neill

Backgrounds in photography brought Colm to digital experiments from the very start of training as a graphic designer. Frustrations with the single and unidirectional workflow of most digital professions, everything to made sense when he's later received as an OSP intern for a few months. Then fully committing to the team, while blancing a student life in all of this. Interests in the digital formats and ways to peel off the layers of the web, he's working on alternative networks to exchange knowledge, networking digitally, and ways to expose their inner workings.

Colophon

Participants :

전혜림 Jeon (Matilda) Hye Rim,
구운희 Gu Wun Hyoi,
김명곤 Kim Myung Gon,
정경빈 Jeong Kyung Bin,
서수호 Seo Soo Ho,
안종민 Ahn Jong Min,
조은지 Jo Eunji,
데이비드 David Dahan,
안유 Ahnu Ahn,
손지선 Ji Sun Son,
루디 Ludivine Loiseau,
피에르 Pierre Huyghebaert,
사라 Sarah Magnan,
콜럼 Colm O'Neill

This publication was printed and
binded on the 4th of October.

You can download a pdf version of
this brochure, all files, how to
manuals and notes, and maybe even
the future productions around it on
[http://osp.constantvzw.org/
workshop/typojanchi-seoul](http://osp.constantvzw.org/workshop/typojanchi-seoul)

Curator :

김영나 Na Kim

Coordinator :

전영주 Emily Young Joo Jeon

All pictures

[http://ospublish.constantvzw.org/
images/TypoJanchi-Seoul-2013](http://ospublish.constantvzw.org/images/TypoJanchi-Seoul-2013)

Thanks to Typojanchi !

Licence

This publication is © by the participants of the workshop unless noted otherwise.

The picture seems to be public domain, or of fair use.

As the copyright holders, we open up these materials for your to use if you want to.

Different types of materials have different licenses. This is to be accomodating to different use cases and cultures that have arisen around different digital formats.

The licenses we use are copyleft. That means that you are free to reuse, modify and redistribute our materials, and that you are also allowed to make money doing so. You have to, however, redistribute your own variants under the same license. This means that the ecology of sharing is stimulated.

All texts written in natural languages, and all drawings and visual designs, are dual licensed under the Free Art License and Creative Commons Share Alike.
- The Free Art License, because we feel that the spirit of its beautiful and simple seems the closest as the way we think creation can be shared.

- Creative Commons, because a lot of people take part in this ecosystem.

If you modify our work, you are free to license your version under either of them (or both).

As our Lafkon state : "This work contains material which may be subject to trademark laws in one or more jurisdictions. Before using this content, please ensure that you have the right to use it under the laws which apply in the circumstances of your intended use."

And they explain it with "We try to keep some way of a sense-making differentiation for different copyrights of different files. this disclaimer is like the standard disclaimer for files that contain things we've not created ourselves (like some logos) or we are not sure whether to put in under our standard licensing. Since we do not want to exclude these files from download, we put the disclaimer."

If you are not certain of the license of a file that you want to re-use, please feel free to ask us directly by email:

<mail@osp.constantvzw.org>.

