

IDSgrep, version 0.4pre

Matthew Skala

March 7, 2013

Contents

Quick start	3
Introduction	5
What's new	5
Download, build, test, and install	5
Unicode IDSes	7
Interface to CHISE IDS	7
Interface to KanjiVG	9
Interface to EDICT2	9
Interface to Tsukurimashou	10
A note on TrueType/OpenType	11
Invoking idsgrep	13
Command-line options	13
Environment variables	14
Technical details	15
The data structure	15
EIDS syntax	15
Matching	18
Match anything	19
Match anywhere	19
Match children in any order	19
NOT	19
AND	19
OR	20
Literal tree matching	20
Associative matching	20
Regular expression matching	21
User-defined matching predicates	22
Cooked output	22
Bibliography	25

Quick start

Use `idsgrep` much as you would use `grep`:

```
idsgrep [<options>] <pattern> [<file> ...]
```

Its general function is to search one or more files for items matching a pattern, like `grep` [7] but with a different pattern syntax. Although potentially usable for an unlimited range of tasks, `idsgrep`'s motivating application is to searching databases of Han script (Chinese, Japanese, etc.) character descriptions. It provides a much more powerful replacement for the “radical search” feature of dictionaries like *Kiten* [8] and *WWWJDIC* [5].

The syntax for matching patterns, and the range of command-line options available, are complicated. Later sections of this document explain those things in detail; for now, here are some examples.

`idsgrep 萌 dictionary`

A literal character searches for the decomposition of that character, exact match only.

`idsgrep -d 萌`

The `-d` option with empty argument searches a default collection of dictionaries.

`idsgrep -dtsuku 萌`

The `-d` option can be given an argument to choose a specific default dictionary. Note the argument must be given in the same `argv`-element with the `-d`; the syntax `-d tsuku` with a space would mean “Use the default dictionaries and search for the (syntactically invalid) pattern ‘tsuku.’”

`othersoft | idsgrep 萌`

Standard input will be used if no other input source is specified.

`idsgrep -d ...日`

Three dots match their argument anywhere, so this will match 日, 早, and 萌.

`idsgrep -d '?'`

A question mark, which will probably require shell quoting, matches anything. This is most useful as part of a more complex pattern.

`idsgrep -d 'ㇿ?心'`

Unicode Ideographic Description Characters can be used to build up sequences that also incorporate the wildcards; this example matches characters consisting of something above 心, such as 忽 and 恋 but not 虑.

`idsgrep -d '[tb](anything)心'`

There are ASCII aliases for operators that may be inconvenient to type; this query is functionally the same as the previous one.

`idsgrep -d '&!萌...|日月'`

Boolean prefix operators include `&` (AND), `|` (OR), and `!` (NOT). This example matches anything that contains 日 or 月 but is not 萌.

`idsgrep -d '*ㇿ日?'`

Asterisk makes children match in any order; this example matches 日 at left or right.

`idsgrep -d '@ㇿㇿ?日?'`

At-sign treats an operator as associative; this example matches both ㇿㇿ?日? and ㇿ?ㇿ日?.

`idsgrep -d '.../(femoral)'`

Slash invokes Perl-compatible regular expression matching, which might be useful for the EDICT2-based meaning dictionary.

`idsgrep -d '...=?'`

Equals escapes matching operators; this example searches for a literal question mark anywhere in the tree.

`idsgrep -d '\X840C'`

Several kinds of backslash escapes allow entering characters that might not otherwise be available.

`idsgrep -d -c indent 萌`

The `-c` option selects “cooked” or pretty-printed output modes.

`idsgrep -d -f FontFile.otf '#1'`

The `-f` option reads the character set of an OpenType font and makes it available as a user-defined matching predicate accessed with the

hash-mark; in the example, it looks up each character in the default dictionaries.

`idsgrep -U '?'`

The `-U` option generates a list of Unicode characters.

`idsgrep -Uxdb '?'`

An optional argument to `-U` specifies information to include in the generated list entries: `x` for hexadecimal value, `d` for decimal, `b` for block name.

`idsgrep -U -f FontFile.otf '#l'`

Combine `-U` and `-f` to list the characters in a font.

Introduction

The Han character set is open-ended. Although a few thousand characters suffice to write the languages most commonly written in Han script languages (namely Chinese and Japanese) most of the time, popular standards define tens of thousands of less-popular characters, and there are at least hundreds of thousands of rare characters known to occur in names, historical contexts, and in languages like Korean and Vietnamese that may still use Han script occasionally despite now being written primarily in other scripts.

Computer text processing systems that use fixed lists of characters will inevitably find themselves unable to represent some text. As a result, there is a need to *describe* characters in a standard way that may have no standard code points of their own. A similar need for descriptions of characters arises when looking up characters in a dictionary; a user may recognize some or all the visual features of a character (such as its parts and the way they are laid out) without knowing how to enter the character as a whole.

IDSgrep's main function is to query character description databases in a flexible way. This need was identified during development of the Tsukurimashou font family [12]; there, the visual appearance of Han character glyphs corresponds directly to the MetaPost code implementing them, and the desire for code re-use and consistency often motivates a close examination of the existing work to answer questions like "What other characters contain this shape, and how did we implement it last time?" Standard tools like `grep` [7] can sometimes be applied to answer such questions by searching for subroutine names in the source code, but the related question of "What other characters, not yet implemented, could we build that would use this shape?" requires comparing with some external database of the characters commonly used in the language. How can we run `grep` on the writing system itself?

Someone confronted with an unknown character and wanting to look it up in a more ordinary dictionary to find the meaning may, similarly, want to search for characters based on specific features while

leaving others unspecified, with questions like "What characters exist that have the common 心 shape at the bottom, with the upper part divided into two things side by side? The two things at the top are shapes I don't recognize, printed too small for me to identify them more precisely." Existing dictionary-query methods are not adequate for some reasonable queries of this nature.

For instance, the radical-and-stroke-count method of traditional character dictionaries requires identifying the head radical and counting strokes, both of which may be difficult; dictionary codes like SKIP and Four Corners key on some layout attributes but not all; multi-radical search allows the user to choose whichever radicals they recognize, but it ignores layout entirely; and computer handwriting recognition generally works well if and only if the user is sure of the writing of the first few strokes in the character. Furthermore, these search schemes often are implemented only in heavy, non-portable, GUI software that cannot be automated and mixes poorly with standard computing environments. IDSgrep can answer the example query correctly with a single, simple command line (`idsgrep -d '[tb][lr]??心'`). This software is intended to bring the user-friendliness of `grep` to Han character dictionaries.

What's new

The main new features in version 0.4 are:

- changes to the build system for better integration with Tsukurimashou;
- support for user-defined matching predicates, and in particular, the ability to match against the list of characters defined by a font file ("f" option);
- built-in generation of Unicode character lists ("-U" option).

Download, build, test, and install

IDSgrep is distributed under the umbrella of the Tsukurimashou project on Sourceforge.JP [12], <http://>

`//tsukurimashou.sourceforge.jp/`. Releases of IDSgrep will appear on the project download page; development versions are available by SVN checkout from the `trunk/idsgrep` subdirectory of the repository. For the convenience of Github users, the Tsukurimashou (and thus IDSgrep) repository is also mirrored into a Github repository [13], but the project on Sourceforge.JP and its SVN repository remain the main public locations for IDSgrep development and all bug-tracker items should be filed there.

A minimal default build and install could run something like this:

```
tar -xvzf idsgrep-0.3.tar.gz
cd idsgrep-0.3
./configure
make
su -c 'make install'
```

IDSgrep can build dictionaries from the Tsukurimashou font package, which is available through the same Sourceforge.JP project as IDSgrep; from the KanjiVG database available at <http://kanjivg.tagaini.net/> [3]; from the CHISE IDS database available at <http://chise.zinbun.kyoto-u.ac.jp/dist/ids/> [1]; or from the EDICT2 database available at <http://www.csse.monash.edu.au/~jwb/edict.html> [4]. For an ideal complete installation of IDSgrep, one would download all those packages, build Tsukurimashou first, and make it and the dictionaries available to the IDSgrep configure script. A precompiled version of the CHISE IDS-derived dictionary is bundled in the IDSgrep distribution tarball, so that one should be available (though not necessarily up-to-date) without any dependencies.

The configure script will by default make a reasonable effort to find the dependencies; in many common cases it is not necessary to specify them on the command line. Here is a more complete installation process relying on configure to find Tsukurimashou in a sibling directory and the others in the current directory:

```
unzip tsukurimashou-0.6.zip
cd tsukurimashou-0.6
./configure
make
# install of Tsukurimashou not needed by IDSgrep
cd ..
tar -xvzf idsgrep-0.2.tar.gz
cd idsgrep-0.2
ln -s /some/where/else/kanjivg-20120219.xml.gz .
```

```
ln -s /some/where/else/edict2.gz .
ln -s /some/where/else/chise-ids-0.25 .
./configure
make
make check
su -c 'make install'
```

It is necessary to at least configure Tsukurimashou, if not fully build it, before building IDSgrep. The IDSgrep build will then invoke the Tsukurimashou build to create just the files needed by IDSgrep. It is not necessary to configure or build CHISE IDS (which would require first installing other parts of the larger CHISE system and probably XEmacs as well); IDSgrep only needs to look at the CHISE IDS data files.

If the default search fails, the filenames of KanjiVG (`.xml` or `.xml.gz`), EDICT2 (`.gz`), and the directories containing extracted distributions of Tsukurimashou and CHISE IDS can be specified on the configure command line with the `--with-kanjivg`, `--with-edict2`, `--with-tsuku-build`, and `--with-chise-ids` options. For other options, run `configure --help`. It's a reasonably standard GNU Autotools [6] configuration script, albeit with a lot of options for inapplicable installation directories removed to simplify the help message.

The EDICT2-based dictionary should preferably include character decompositions from some other dictionary; which one is selectable by the `--enable-edict-decomp` option. Allowed values include `chise`, `kanjivg`, `tsuku`, and `no`; the default of `auto` will try all of those in that order and use the first that works. The value `no` corresponds to simply mapping every character to itself without further decomposition; that is obviously not as informative as might be desired, but it will still allow for regular expression searches.


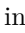
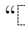


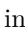

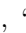
The “check” Makefile target runs the IDSgrep test suite. Some tests require the dictionary files and will be skipped if those are not present. There is also a test that will use Valgrind [11] if available, to check for memory-related problems; if Valgrind is not found in the PATH, this test will be skipped.

The configure script supports an `--enable-gcov` switch to enable meta-testing of the test suite's coverage. This feature requires that the Gcov coverage analyser be installed. To do a coverage analysis, run configure with `--enable-gcov` and any other desired options, then do `make clean` (necessary to be sure all object files are rebuilt with the coverage instrumentation) followed by `make check`. Full coverage can only

be attained if the dictionary files are installed (not just built). Most people would not want to install the IDSgrep binary itself when built under this option.

Unicode IDSes

Although IDSgrep uses a more elaborate syntax, it is well to know about the Unicode Consortium’s “Ideographic Description Sequences” (IDSes), which are a subset of IDSgrep’s. These are documented more fully in the Unicode standard [14].

- A character from one of the Unified Han or CJK Radical ranges is a complete IDS and simply represents itself. For instance, “大” is a complete IDS.
- The Ideographic Description Characer (IDC) code points U+2FF0, U+2FF1, and U+2FF4 through U+2FFB, whose graphic images look like , are prefix binary operators. One of these characters followed by two complete IDSes forms another complete IDS, representing a character formed by joining the two smaller characters in a way suggested by the name and graphical image of the IDC. For instance, “日月” describes the character 明. These structures may be nested; for instance, “言五口” describes the character 語.
- The IDC code points U+2FF2 and U+2FF3, which look like , are prefix ternary operators. (Unicode uses the less-standard word “trinary” to describe them.) One of them can be followed by three complete IDSes to form an IDS that describes a character made of three parts, much in the same manner as the binary operators. For instance, “言糸糸久” describes the character 變.
- As of Unicode 6.1, IDS length is unlimited. Earlier versions specified that an IDS could not be more than 16 code points long overall nor contain more than six consecutive non-operator characters. This rule appears to be have been intended to make things easier for systems that need to be able to jump into the middle of text and quickly find the starts and ends of IDSes.
- IDSes non-bindingly “should” be as short as possible and should reflect “the natural radical-phonetic division for an ideograph if it has one.”

The Unicode standard does not mention variation selectors in any IDS-related context, except that it offers the possibility of prefixing U+303E, the “ideographic variation mark,” to the entire sequence to indicate a variation. Such a prefix is explicitly defined not to be counted as part of the IDS.

My opinion is that Unicode did not intend to permit variation selectors inside IDS syntax. Variation selectors arguably exist to patch over encoding inadequacies resulting from Unicode’s internal politics. When a code point is not really specific enough, because it refers to two or more things which you think are not actually the same thing, then you can add a variation selector to indicate which thing you really mean. IDSes, on the other hand, bearing in mind that they are imported from GBK, operate at a lower level to specify characters in terms of parts that are assumed to be adequately encoded. If a code point to be used in an IDS is not specific enough, then that element should be described with a smaller fragment of IDS syntax instead of by using the ambiguous code point. If the closest match possible is still not perfect, then it is time to use U+303E. The fact they offer the U+303E mechanism for specifying variations offers further support to the idea that they did not intend to allow variation selectors *inside* IDSes.

However, it’s a difficult question because IDSes, by addressing the visual appearance of characters instead of their semantics, fundamentally challenge the basic Unicode principle that code points specify characters and not glyphs. The distinction between characters and glyphs simply cannot be made perfectly in all cases. For use in cases where variation selectors appear to be appropriate, both CHISE IDS and IDSgrep extend IDS syntax in such a way as to allow them in some way.

Interface to CHISE IDS

The CHISE project [1] maintains a database of Han characters covering multiple languages as part of a larger processing environment that also includes a version of XEmacs [15] modified to follow the principles of the UTF-2000 initiative [9]. It also has connections to GlyphWiki [2]. These systems are documented primarily in Japanese; English-language documentation is sparse and not necessarily up to date.

For IDSgrep’s purposes, the most interesting part of CHISE is a module called CHISE IDS, which includes a database of about 140000 characters (exact count depending on the version), with decompositions in its own extension of Unicode IDS syntax.

The main purpose of this IDS database is to provide a search capability within the modified XEmacs; there is also code for a Web search form. From examination of the database files it appears that the rules for CHISE IDS's extended IDS syntax are more or less as follows.

- Generally, Unicode IDS rules apply.
- An XML entity-like sequence of the form `&NAME;` counts as a single ideograph. The field indicated by `NAME` is a symbolic identifier or database key defined internally by the project. Such identifiers have been observed to contain uppercase ASCII letters, numerals, hyphens, and plus signs; they usually consist of a short alphabetic prefix, a hyphen, and a number. These entity references are usually used to refer to characters for which CHISE has an encoding and Unicode doesn't.
- A Unicode variation sequence (an ideograph followed by a variation selector) counts as a single character.

Although CHISE IDS's extensions to IDS permit strings that would not be valid IDSgrep EIDS syntax, it is easy to convert them into EIDS format. IDSgrep includes a `chise2eids` Perl script for that purpose. The configure script will look for CHISE IDS in a directory named `chise-ids-*` in a short list of likely places, or use the value of the `--enable-chise-ids` command-line option if one is given. This directory should simply be an unpacked CHISE IDS distribution tarball, or a checkout from the CHISE IDS Git repository. It is not necessary to run CHISE's Makefile, which would require also having and installing other parts of the larger system.

As of this writing (March 6, 2013), the CHISE IDS distribution tarballs are available from <http://chise.zinbun.kyoto-u.ac.jp/dist/ids/>, and there is a Git repository at <http://git.chise.org/git/chise/ids.git>. However, on March 5, there was an announcement posted to the CHISE Project's mailing lists in English and Japanese saying that on March 11 all the servers at chise.org (which includes the mailing lists themselves) will cease to operate. It is not clear to me whether that is intended to be something temporary, or the permanent end of the project, but given that nothing seems to have happened in CHISE in years, I fear the worst. I have posted a snapshot of the current CHISE IDS Git repository in my own Github account at <https://github.com/mskala/chise-ids>;

that should be a long-term stable source of the data used by IDSgrep.

The last distribution tarball of CHISE-IDS was version 0.25 dated June 2010. The Git versions are more recent and may be preferable. The directory created by checking out a Git version will probably not have a name recognized automatically by the build system, so it should be given on the configure command line with `--enable-chise-ids`.

Roughly 6% of the entries in the CHISE IDS database include invalid extended IDS syntax—most often in the form of too many children for the operators used, or less often, too few. Most but not all of the errors occur in the `IDS-HZK??.txt` files, which are unmaintained. It appears that the native search tools for the database generally work on the basis of pure substring searches, where the higher-level syntax errors that would be detected by the IDSgrep parser can go unnoticed. The `chise2eids` program generates a `chise.errs` file during build, listing all the syntax errors it finds (11748 of them in the current Git version as of this writing); invalid entries are otherwise ignored and will not appear in the main output file `chise.eids`. Although 6% may sound like a lot of errors, the invalid entries are generally in sufficiently obscure character components that it should have little practical effect on the quality of dictionary lookups: at worst, some character components may end up not broken down into pieces as small as would otherwise be possible.

CHISE IDS refers to individual characters in a more general way than just by single Unicode code point: sometimes it uses a variation sequence consisting of a kanji code point followed by a variation selector in the `U+FE00` to `U+FE0F` or `U+E0100` to `U+E01EF` ranges, and sometimes it uses a string that looks like an XML character entity reference, along the lines of `&NAME;`. Both of these map naturally to IDSgrep's concept of a multi-character head. The two-code-point sequence `U+840C U+E0101` is translated to the IDSgrep syntax `<X840cX{E0101}>;`, and the XML-like syntax `&FU-123;` is translated to `<FU-123>;`. CHISE IDS does not seem to refer to the same character in different ways (for instance, a code point with no variation selector somehow matching as a default to the same code point with a variation selector, which might be plausible under Unicode's definition of what variation selectors signify) and `chise2eids` does not attempt to accomodate anything like that.

The CHISE IDS database is covered by the GNU

GPL version 2 or later, which is basically compatible with the GNU GPL version 3 used by IDSgrep. So distributions of IDSgrep can reasonably include a bundled copy of `chise.eids` for the benefit of users who don't want to download the separate package and generate their own. Without taking a position on whether the `chise.eids` file constitutes "object code" for the purposes of the GNU GPL as opposed to being modified source code in itself, I am willing to provide a copy of the CHISE IDS checkout I used to generate my version of `chise.eids`, to anyone who contacts me about it at `mskala@ansuz.sooke.bc.ca`. Going directly to the original distribution points at the URLs given above is probably a more convenient option for most users.

Interface to KanjiVG

The KanjiVG project [3] maintains a database of kanji (Han characters as used by Japanese) in an extended SVG format, which implies that it is XML. The `kvg2eids` Perl script, included as part of IDSgrep, is capable of reading this database and converting it to Extended Ideographic Description Sequences (EIDSeS). As described above, if a reasonably recent version of KanjiVG's compressed XML file is available to configure, then IDSgrep's build will create such a dictionary and `make install` will install it.

KanjiVG describes characters primarily in terms of strokes, not multi-stroke components, and it attempts to follow the official stroke order and etymological component breakdown. That approach results in some peculiarities from the point of view of dictionary searching. For instance, in the kanji 囧, the official stroke order is to write two strokes of the enclosing box, then the central glyph, then the bottom of the box. KanjiVG's XML file lists two "elements" identified with the kanji 口, one for the first two strokes and one for the final stroke, with additional attributes specifying that they are actually two parts of the same element. KanjiVG has changed its own standard for how to represent this information in the recent past, and not all entries have been updated to the latest standard yet. The current version of `kvg2eids` does not correctly process 囧 nor some other characters with parts written in nonsequential order. On that particular one it generates a special functor containing debugging information; for some others, it may actually generate an EIDS with the same radical appearing multiple times, following the structure described in KanjiVG whether it's what was intended or not. As a result, not all entries in the dictionary

will be right. However, only a few are affected by this issue.

As of March 2012, I (Matthew Skala, the author of IDSgrep) have become a member of the KanjiVG project and there is some possibility that KanjiVG's database design will change in a way that makes it easier to recover spatial organization for searching with IDSgrep.

With the current versions of IDSgrep and KanjiVG, the KanjiVG-derived dictionary contains 6660 entries covering all the popularly-used Japanese kanji. Note that the KanjiVG input file, and presumably the resulting format-converted dictionary, are covered by a Creative Commons Attribution-ShareAlike license, distinct from the GNU GPL applicable to IDSgrep itself.

Interface to EDICT2

Jim Breen's JMdict/EDICT project maintains a file called EDICT2 [4] which is more like a traditional dictionary, with words and meanings, than a database of kanji. Such dictionaries are not the primary target of IDSgrep and IDSgrep's query syntax is not perfectly suited to them. However, the regular-expression matching features may make it practical to search EDICT2 with IDSgrep, and there is some value in being able to do sub-character structural searches on multi-character words.

If another dictionary besides EDICT2 is available (subject to configuration by `--enable-edict-decomp`), then the build system will generate and install a dictionary file called `edict.eids` which represents a database join of EDICT2 with the other dictionary. With no other dictionary, the file can still be generated but will contain no character decomposition information. A sample entry might look like this:

【明】.<明>日月 ([みん] (n) Ming (dynasty of China))

The head for the entire entry is the head from the EDICT2 entry. Then the tree is a binary tree with a comma as the functor and the first child being the entire decomposition dictionary entry for the first character. The second child represents the rest of the entry. With a two-character or longer head, this child would also be a binary comma with the second character of the entry head as its first child. In this way the characters of the entry head are all represented as left children of commas, forming a linked-list structure (much like a Prolog linked-list with commas instead of dots as the functors). The final child at the

bottom is a nullary node containing as its functor simply the rest of the EDICT2 entry.

The rationale for this syntax is that it allows a relatively simple way of querying multi-character words in EDICT2 using the existing IDSgrep query types. To find an exact match, just query the head (which will require head brackets and a semicolon if the query is more than one character long), as in `idsgrep -de '<教育>;'`. To search for the first few characters, commas can be imagined as separators (though their actual function is quite different) with a comma at the start and a question mark at the end, as in `idsgrep -de ',教育?'`. These queries can be combined with the sub-character breakdown queries already supported by the decomposition dictionaries. For instance, `idsgrep -de ',教,...|日月!??'` will search for, and give definitions of, words of exactly two characters in which the first is 教 and the second character contains 日 or 月 anywhere. The restriction to exactly two characters is accomplished by the sub-query “!??”, which fails to match on the binary comma that would be present at that point in a longer word.

EDICT2 is under the Creative Commons Attribution-ShareAlike license. Since KanjiVG is as well, that license would presumably also apply to a combined dictionary made from EDICT2 and KanjiVG. An EDICT2-only dictionary with no decompositions from other sources should similarly be under Creative Commons Attribution-ShareAlike. It might not be legal to distribute outside one’s own organization a dictionary formed by joining EDICT2 with CHISE IDS or Tsukurimashou, because those sources are covered by versions of the GNU GPL, which is not compatible with the Creative Commons license.

Interface to Tsukurimashou

IDSgrep is closely connected with the Tsukurimashou font family [12]. They have the same author; it was largely for use in Tsukurimashou development that IDSgrep was developed at all; and IDSgrep’s source control system is a subdirectory within Tsukurimashou’s. Building IDSgrep in conjunction with Tsukurimashou allows IDSgrep to extract from the Tsukurimashou build system a dictionary of character decompositions as they appear in Tsukurimashou. The Tsukurimashou fonts are also necessary to build this IDSgrep user manual. However, IDSgrep is also distributed as a separate package, because it will be of use to non-users of Tsukurimashou, and the Tsukurimashou build system will not recurse into IDSgrep’s

directory and build IDSgrep by default; only if requested.

IDSgrep is one of several parasite packages of Tsukurimashou, using a mechanism introduced in Tsukurimashou 0.7 and IDSgrep 0.4. Previous versions used a different interface.

To build Tsukurimashou with IDSgrep: specify the “--enable-parasites” option to Tsukurimashou’s configure script with an appropriate value, such as “--enable-parasites=idsgrep”. See the Tsukurimashou documentation for other possible values of this option. Building Tsukurimashou will then implicitly build IDSgrep. It should be possible to pass IDSgrep configure options to Tsukurimashou’s configure script and have them automatically passed down the chain (in the standard Autotools sub-package fashion) but that is not well-tested.

To build Tsukurimashou without IDSgrep: this is the default when you run the Tsukurimashou build from the root of the Tsukurimashou distribution. The IDSgrep source is included as a subdirectory in distributions of Tsukurimashou, but only built on request.

For a more customized build of IDSgrep, with or without Tsukurimashou: you can also run IDSgrep’s configure in its own directory, and then do make (and the usual targets) there. It will look for Tsukurimashou (specifically, a build directory in which Tsukurimashou’s configure has *already been executed*) as the parent directory and in a few other places, or you can specify the location of a Tsukurimashou build with the “--with-tsuku-build” option to IDSgrep’s configure. If Tsukurimashou is not available, IDSgrep will build without creating the Tsukurimashou-derived dictionary file.

During IDSgrep’s build, if it can access a Tsukurimashou build directory, it will recursively call make on Tsukurimashou’s build system. That is a hook that causes Tsukurimashou’s build system to generate the EIDS decomposition dictionary, which is then copied or linked back into IDSgrep’s build directory and can be installed with IDSgrep’s `make install`. IDSgrep’s build will also look in Tsukurimashou’s build directory for the font “Tsukurimashou Mincho” which is needed to build this user manual, and will make recursive calls to make for Tsukurimashou to build that if necessary. This kind of upward-callback make invocation is a little inefficient (in particular, it does not handle jobserver mode well) so it is better, if you want both packages, to use the centralized Tsukurimashou build system, which will do its own thing first and then call IDSgrep’s build near the end in a better-

integrated way. If you want to run “make install” just on IDSgrep and not on Tsukurimashou (which might be a reasonable thing to want because of operating system font installation issues), you should run just “make” in Tsukurimashou’s directory, then cd to IDSgrep’s directory and run “make install.”

A note on TrueType/OpenType

This version of IDSgrep is designed to read TrueType or OpenType files (the distinction between the two is not relevant at this level) for character map information. The specification for the TrueType/OpenType file format reads like a parody. I’d like to take a moment to complain about a few things.

- Although the format contains binary fields which must be read in a specific byte order, one of the two magic numbers that can identify the file format (for a single font as opposed to a “collection”) is 0x4F54544F, which is a palindrome at the byte level and thus useless for detecting byte order problems.
- The other possible magic number for non-collection files is 0x00010000, which is quite likely to occur in files that are not TrueType/OpenType files, making it harder to detect when one may have been passed a bad file.
- Many decades of research on error detection codes were ignored in the design of the OpenType checksum algorithm, which (among other issues) cannot detect any reordering of 32-bit words unless it crosses a table boundary. At least the algorithm produces its meaningless results fast; yay, efficiency!
- There are 32-bit byte offsets referenced to the start of the file. There are 32- and 16-bit byte offsets referenced to the start of the current table. There are 32- and 16-bit byte offsets referenced to the *locations of the offset fields themselves*, so a field at offset 0x1234 referring to another field at offset 0x5678 will contain 0x4444. There are also indices measured in units other than bytes.
- There are variable-length objects not tagged with their lengths except indirectly: they are presumably contained entirely within larger objects that are tagged with lengths.
- Consider the cmap format 4 subtable, which Microsoft says is their preferred format. It includes

four variable-length arrays each containing segCount number of two-byte entries. The value of segCount is not directly recorded anywhere, but these values are all required in the header:

- $2 \cdot \text{segCount}$;
- $2 \cdot 2^{\lfloor \log_2 \text{segCount} \rfloor}$;
- $\log_2(2 \cdot 2^{\lfloor \log_2 \text{segCount} \rfloor} / 2)$ (which is described like that in the spec); and of course
- $2 \cdot \text{segCount} - 2 \cdot 2^{\lfloor \log_2 \text{segCount} \rfloor}$.

- The bizarre length-derived values in the format 4 header (and other similar sets of table-size-logarithm numbers that occur elsewhere in the file format) appear to be designed to support someone’s binary search code. Instead of computing those numbers itself starting from the length, the search code can just use values straight from the table to initialize its variables. Consider what would happen if someone actually did that as the designers apparently intended, and the numbers happened to be incorrect in the file. If, for instance, numbers in the file were swapped around on 32-bit boundaries, the checksums wouldn’t detect a problem; and the speed demons who think they need precomputed logarithms probably aren’t wasting time checking checksums anyway. The code isn’t checking whether the numbers are consistent (because to do that you would have to calculate them fresh, and then why bother storing them in the first place?), so it will end up “searching” into random areas of the file, or into uninitialized memory beyond. Now think about the relative costs of disk reads, network transfer, and arithmetic, and consider whether having those values precalculated and stored in the file would actually save any time even if they could be trusted.
- The cmap format 4 subtable consists of, in this order: fixed-length stuff totalling 14 bytes; one variable-length array of length $2 \cdot \text{segCount}$ bytes; *one more two-byte fixed-length field*; three more variable-length arrays each of length $2 \cdot \text{segCount}$ bytes; and finally, one more variable-length array whose length is not directly specified anywhere but could presumably be inferred by subtracting from the known size of the overall table. The four $2 \cdot \text{segCount}$ -byte arrays are actually the rearranged slices of a single logical array whose elements are four-field structures; but the extra

reserved two bytes stuck in the middle of the table make a straightforward transposition impossible. Four-tuples of the same kind with the same four fields also occur in the format 2 sub-table; but there, they occur as a single array with each record written in an 8-byte block.

- It is an intended, documented feature that some of the variable-length arrays in TrueType/OpenType may overlap with each other. As a result, bounds-checking, in addition to being intrinsically difficult because of the lack of information, would cause the reader to reject some files that the specification claims are legitimate.
- Code-injection bugs allowing execution of arbitrary code in a privileged context have been reported in software that implemented this file format without bounds-checking. This should surprise no one.

IDSgrep attempts to do all reasonable bounds-checking on the fields it needs, and to ignore fields it does not need; given a bad TrueType/OpenType file, it is intended that IDSgrep should be able to make the best of it and at worst fail gracefully with an error message. It should not be possible to crash IDSgrep by giving it a bad font file to read.

However, the nature of the file format means that at least in the current version, we can't be confident all possible problems have been foreseen and excluded. Let me know if you find a font file that makes IDSgrep crash and I'll try to fix it. IDSgrep probably should not be allowed to read font files supplied by untrusted sources such as Web users.

Invoking idsgrep

The command-line `idsgrep` utility works much like most other command-line programs, and like `grep` [7] in particular. It takes options and other arguments. The first non-option argument is an EIDS representing the matching pattern, and any remaining non-option arguments are taken as filenames to read. If there are no filenames, `idsgrep` will read from standard input. Output always goes to standard output.

When there is more than one file being read (either by direct specification or indirectly with the `-d` dictionary option), `idsgrep` will preface each EIDS in its output with “:*<filename>*:” to indicate in which file the EIDS was found. Note that under the EIDS syntax rules, that creates a unary node senior to the entire tree, so that the output remains in valid EIDS format, except in the case of filenames containing colons, which will be handled via backslash escapes in the future when those are fully implemented for output.

Command-line options

`-d, --dictionary` Read a dictionary from the standard location. There is a pathname for dictionaries hardcoded into the `idsgrep` binary, generally `{prefix}/share/dict`, and if this option is given, its argument (which may be empty) will be appended to the dictionary directory path, followed by “`*.eids`,” and then treated as a shell glob pattern. Any matching files are then searched in addition to those otherwise specified on the command line. A small added wrinkle is that when more than one file is searched (resulting in `:filename:` tags on the output lines), any of them that came from the `-d` option will be abbreviated by omitting the hardcoded path name. The purpose of this option is to cover the common case of searching the installed dictionaries. Just specifying “`-d`” will search all the installed dictionaries; specifying an abbreviation of the dictionary name, as “`-dt`” or “`-dk`,” will search just the matching one; and it remains possible to specify a file exactly or use standard input in

the usual `grep`-like way.

`-c, --cooking` Select the output generation and input canonicalization mode. Requires one argument, which may be one of the keywords `raw`, `rawnc`, `ascii`, `cooked`, or `indent`, to specify a preset mode; or a string of up to twelve decimal digits to control the output in more detail. The default mode is `raw`. See the section on “cooked output” in this manual for more details.

`-f, --font-chars` Read a font file and make its character coverage available as a user-defined matching predicate through the “`#`” matching operator. In the current version, this feature can only read TrueType and OpenType files that contain Unicode (or near equivalent) mappings described with cmap subtable types 0, 2, 4, 12, or 13. This option may be specified multiple times, with successive invocations corresponding to user-defined predicates 1, 2, 3, and so on. The maximum number of user-defined predicates is limited to the number of bits in the largest integer type available to the C compiler; 32 or 64 on many systems.

`-U, --unicode-list` Generate a dictionary of Unicode code points, and read that before reading any other dictionaries or input files that may be specified. The generated dictionary consists of a single line for each of the code points U+0000 through U+10FFFF in ascending order, excluding the surrogates but not any other invalid or non-character code points; on each line, there is a tree whose head is the character and whose body is either a nullary semicolon or (if the optional argument to `-U` was specified) a nullary functor containing semicolon-separated pieces of information selected by the characters of the optional argument. Characters permitted in the argument are “`b`” for the Unicode block name; “`d`” for the decimal value of the code point; and “`x`” for the hexadecimal value with “`U+.`” For example, specifying “`-Uxdb`” will generate and scan a dic-

tionary that includes the line “<A>(U•0041;65;Basic latin).” This option is intended to be used together with -f to produce font coverage lists.

-V, --version Display the version and license information for IDSgrep.

-h, --help Display a short summary of these options.

Environment variables _____

The idsgrep utility recognizes just one environment variable, IDSGREP_DICDIR, which if present specifies a directory for the -d option to search instead of its hardcoded default.

Note that idsgrep does not pay attention to any other environment variables, and in particular, not LC_ALL and company. The input and output of this program are always UTF-8 encoded Unicode *regardless of locale settings*. Since the basic function of this program is closely tied to the Unicode-specific “ideographic description characters,” it would be difficult if not impossible for it to work in any non-Unicode locale. Predictability is also important because of the likely usefulness of this software in automated contexts; if it followed locale environment variables, many users would have to carefully override those all the time to be sure of portability. Instead of creating that situation, idsgrep by design has a consistent input and output format on all systems and users are welcome to pipe things through a conversion program if necessary.

Technical details

This section is intended to describe IDSgrep’s syntax and matching procedure in complete precise detail; and those things are, in turn, designed to be powerful rather than easy. As a result, the description may be confusing for some users. See the examples in the “Quick start” section for a more accessible introduction to how to use the utility.

The system is best understood in terms of three interconnected major concepts:

- an abstract data structure;
- a syntax for expressing instances of the data structure as “Extended Ideographic Description Sequences” (EIDSes);
- a function for determining whether two instances of the data structure “match.”

Then the basic function of `idsgrep` is to take one EIDS as a matching pattern, scan a file containing many more, and write out the ones that match the matching pattern. The three major concepts are described, one each, in the following sections. A final section describes options for how the command-line `idsgrep` program generates EIDS syntax on output.

The data structure

An *EIDS tree* consists of the following:

- An optional *head*, which if present consists of a nonempty string of Unicode characters.
- A required *functor*, which is a nonempty string of Unicode characters.
- A required *arity*, which is an integer from 0 to 3 inclusive.
- A sequence of *children*, of length equal to the arity (no children if arity is zero). Each child is, recursively, an EIDS tree.

Trees with arity zero, one, two, and three are respectively called nullary, unary, binary, and ternary.

Note that these “nonempty strings of Unicode characters” will very often tend to be of length one (single characters) but that is not a requirement. They cannot be empty (length zero); the case of a tree without a head is properly described by “there is no head,” not by “the head is the empty string.” *At present* no Unicode canonicalization is performed, that being left to the user, but this may change in the future. Zero bytes (U+0000) are in principle permitted to occur in EIDS trees, but because Unix passes command-line arguments as null-terminated C strings, they can only be entered in matching patterns via backslash escape sequences.

Typically, these trees are used to describe kanji characters. The literal Unicode character being described will be the head, if there is a code point for it; the functor will be either an ideographic description character like `□` if the character can be subdivided, or else nullary `;` if not. Then the children will correspond to the parts into which it can be decomposed. Some parts of the character may also be available as characters with Unicode code points in their own right; in that case, they will have heads of their own.

EIDS syntax

Unicode’s IDS syntax serves a similar purpose to IDSgrep’s extended IDS syntax, but it lacks sufficient expressive power to cover some of IDSgrep’s needs. Nonetheless, EIDS syntax is noticeably derived from that of Unicode IDSes. Broadly speaking, EIDSes are IDSes extended to include heads (which we need for partial-character lookup); bracketed strings as functors (which we need for capturing arbitrary data); and with arbitrary limits on allowed characters and length relaxed (needed for complex characters and so that matching patterns can be expressed in the same syntax).

Here are some sample EIDSes:

大
□田□虫□土土
□厂□今止
【萌】□+<明>□日月
【店】□广<占>□卜口

```

☐+☐☐日?
&...男...女
[tb]+[or][lr]?日[lr]日?

```

The first three of these examples are valid in the Unicode IDS syntax. The next two contain heads, and are typical of what might exist in a dictionary designed to be searched by the `idsgrep` command-line utility. The last three might be matching patterns a user would enter.

EIDS trees are written in a simple prefix notation that could be called “Polish notation” inasmuch as it is the reverse of “reverse Polish notation.” To write a tree, simply write the head if there is one, the functor, and then if the tree is not nullary, write each of the children. Heads and the functors of trees of different arity are (unless otherwise specified below) written enclosed in different kinds of brackets that indicate the difference between heads and functors, and the arity of the tree when writing a functor.

The basic ASCII brackets for heads and functors are as follows:

head	<	>	<example>
nullary functor (0)	()	(example)
unary functor (1)	.	.	.example.
binary functor (2)	[]	[example]
ternary functor (3)	{	}	{example}

Note that the opening and closing brackets for unary functors are both equal to the ASCII period, U+002E.

Some sequences of Unicode characters beginning with “\” (ASCII backslash, U+005C) are treated specially. Backslash followed by a character from a short list of ASCII Latin letters introduces an escape sequence used to substitute for a character that would otherwise be hard to type; backslash followed by any other character (including a second backslash) is equivalent to the other character, but without any special meaning it would otherwise have had. Thus, backslash can be used for instance to include literally in a bracketed string the closing bracket that otherwise would mark the end of the string.

The backslash-letter escapes are listed below. Note that the letters identifying the type of escape sequence are case-sensitive, and all are lower-case except “\X.” However, for sequences that take a parameter, the parameters are not case-sensitive. Note that all characters inside an escape sequence must be literal ASCII, except in the “default” case of a single backslash used to escape a single non-ASCII character. It is not permitted to use recursive backslash escapes to create some of the characters that make

up a multi-character escape sequence like “\x{f}.”

\a	ASCII BEL (U+0007)
\b	ASCII BS (U+0008)
\cX	ASCII control character X
\e	ASCII ESC (U+001B)
\f	ASCII FF (U+000C)
\t	ASCII HT (U+0009)
\n	ASCII LF (U+000A)
\r	ASCII CR (U+000D)
\xHH	two-digit Unicode hex
\XHHHH	four-digit Unicode hex
\x{Hx} \X{Hx}	variable-length Unicode hex

The `\c` escape takes a parameter consisting of a single ASCII Latin letter character (only); it is equivalent to typing Ctrl plus that letter (case insensitive) on a standard keyboard, that is the ASCII control code in the range U+0001 to U+001A obtained by subtracting 64 from the uppercase letter’s ASCII code or 96 from the lowercase letter’s ASCII code.

The hexadecimal escapes `\x` and `\X` offer a choice of two-digit, four-digit, or variable-length (enclosed by curly braces) hexadecimal specification of Unicode code points. The hex codes are case-insensitive. Values greater than 10FFFF, and therefore outside the Unicode range, will be replaced by the Unicode replacement character U+FFFD.

Parsing of bracketed strings has a few features worth noting. First, there is no special treatment of nested brackets. After the “<” that begins a head, for instance, the next unescaped “>” will end the head, regardless of how many other instances of “<” have been seen. However, because no head or functor can be less than one character long, a closing bracket immediately after the opening bracket (which would otherwise create an illegal empty string) is specially treated as the first character of the string and *not* as a closing bracket. Thus, “(0)” is legal syntax for a functor equal to a closing parenthesis, in a nullary tree; and “...” is a functor equal to a single ASCII period in a unary tree, an important example because it is the commonly-used match-anywhere operator. A bracket character specified via a backslash escape, whether by preceding the literal character with a backslash or by giving its hexadecimal code in a “\x” or “\X” construction, is never taken to start or end a bracketed string.

Each pair of ASCII brackets also has two pairs of generally non-ASCII synonyms, as follows:

have aliases. If a functor and arity matches one of the aliases on the following list, it will be replaced with the indicated single-character functor. The idea is to provide verbose ASCII names for single-character functors of special importance to the matching algorithm. Note that the single-character versions are always the canonical ones, and although the brackets are shown explicitly for clarity, they are nearly all characters from the “sugary implicit” list. This feature may be disabled or modified using some settings of the “-c” command-line option; see the section on output cooking for more information.

(anything)	⇒	(?)	.anywhere.	⇒	...
.not.	⇒	!.	.regex.	⇒	./.
.equal.	⇒	.=.	.unord.	⇒	.*.
.assoc.	⇒	.@.	.user.	⇒	.#.
[and]	⇒	[&]	[or]	⇒	[]
[lr]	⇒	[[]]	[tb]	⇒	[[]]
[enclose]	⇒	[[]]	[wrapu]	⇒	[[]]
[wrapd]	⇒	[[]]	[wrapl]	⇒	[[]]
[wrapul]	⇒	[[]]	[wrapur]	⇒	[[]]
[wrapll]	⇒	[[]]	[overlap]	⇒	[[]]
{lcr}	⇒	{[]}	{tcb}	⇒	{[]}

The `idsgrep` command-line utility attempts to follow Postel’s Law with respect to byte sequences that are not valid UTF-8: “be conservative in what you do, be liberal in what you accept from others.” [10] Jesus of Nazareth stated a similar principle somewhat earlier.* Accordingly, invalid UTF-8 on input is not in general treated as a fatal error. Handling of invalid UTF-8 represents a delicate balance of security issues: if invalid UTF-8 is treated as completely fatal, that creates the possibility for denial of service attacks, but if it is permitted to too great an extent, it can create opportunities for things like buffer overflows. In general, the `idsgrep` utility will not itself break when given bad UTF-8, nor will it make matters worse compared to a system that did not include `idsgrep`, but `idsgrep` cannot be counted on to actively protect some other piece of software that would otherwise be vulnerable to bad UTF-8.†

The parser will skip over (as if they did not exist at all) byte sequences that are not valid UTF-8, including the forbidden bytes 0xC0, 0xC1, and 0xF5 through 0xFF; continuation bytes outside valid multibyte sequences; “overlong” sequences (those that would otherwise be valid, but encode a given

code point other than in the shortest possible way); surrogates; and sequences that encode code points outside the Unicode range. Depending on where they occur within a multibyte sequence, some of these things may result in the whole sequence being skipped instead of just the bad bytes, with the parser making its best guess as to what that means. Be aware that some other software may treat some of these things as valid.

When a code point outside the Unicode range, or a surrogate, is specified using a backslash hexadecimal escape, the parser will interpret it as if the substitute character U+FFFD had been specified instead. All UTF-8 sequences *actually generated by* the `idsgrep` program are guaranteed to be valid UTF-8, barring serious programming errors; and matching operations including PCRE matches occur only on the parsed internal representation which is valid UTF-8. Note that PCRE, despite having a deprecated syntax for sub-encoding byte matching, *cannot* be used to detect invalid bytes that the `idsgrep` parser skipped; it sees only what the parser validly parsed. However, since in its default mode the `idsgrep` program will echo through to the output the exact input byte sequence that was parsed to create a tree, not the internal representation, it is possible that non-UTF-8 input could result in non-UTF-8 output. Several cooked output modes, in which `idsgrep` generates its own UTF-8 from the internal representation and provides guarantees of valid UTF-8 or even valid ASCII output, are available but non-default.

Some byte sequences that are valid UTF-8 but not valid Unicode, for instance the sequence that encodes a reversed byte order mark, may possibly go undetected in the input and be allowed in the output, even when cooked, by the current version of `idsgrep`. It is intended that `idsgrep` should detect that kind of thing where it is reasonable to do so, and future versions may do it better than this one does; but some higher-level errors in Unicode usage, such as misuse of combining characters or variation selectors, will probably never fall within the scope of `idsgrep`.

Matching

The basic function of the `idsgrep` command-line utility is to evaluate each item in the database against a matching pattern. The matching patterns are similar in spirit to the “regular expressions” common throughout the Unix world; however, for theoretical and practical reasons standard regular expressions would be unsuitable for the applications considered

*“There is nothing from without a man, that entering into him can defile him: but the things which come out of him, those are they that defile the man.” (Mark 7:15, KJV)

†Genesis 4:9.

by IDSgrep.

The main theoretical issue is that IDSes, whether IDSgrep-style “extended” or Unicode-style traditional ones, belong to the class of *context-free* languages. They describe tree-like structures nested to arbitrary depth, similar in nature to programming-language expressions containing balanced parentheses although balanced parantheses as such are not actually part of EIDS syntax. The natural way to parse these involves an abstract machine with a stack-like memory that can assume an infinite number of different states. Regular expressions can only be used to recognize the smaller, simpler class of *regular* languages, parsable by an abstract machine with a finite-state memory. It is not possible to write a correct regular expression that will match balanced parentheses. Some advanced software implementations of so-called “regular expressions” (for instance, Perl’s) contain special features that make them more powerful than the standard theoretical model, so that they are capable of recognizing some languages that are non-regular, including balanced parentheses. It is also possible to fake a stack with a finite depth limit by writing a complicated regular expression, and that may be good enough in some practical cases. Some users may also settle for just doing a substring query with `grep` and calling the result close enough. But IDSgrep tries to do it in a way that is really right, and that is described precisely in this section.

We will define a function $match(x, y)$ which takes two EIDS trees as input and returns a Boolean value of true or false. We call x the *pattern* or *needle* and y the *subject* or *haystack*. The `idsgrep` command-line utility generally takes x from its command line and repeatedly evaluates this function for each EIDS it reads from its input; it then writes out all the values of y for which $match(x, y)$ is true.

The $match(x, y)$ function is defined as follows:

- If x and y both have heads, then $match(x, y)$ is true if and only if their heads are identical. Nothing else is examined (in particular, not the children). Then the two cases below do not apply.
- If x and y do not both have heads, then $match(x, y) = match'(x, y)$, whose value generally depends on the functor and arity of x . The $match'$ function has many special cases described in the subsections below, expressing different kinds of special matching operations. These operations roughly correspond to the ASCII char-

acters with sugary implicit brackets in EIDS syntax. They are shown with brackets for clarity in the discussion below, but users would generally type them without the brackets and depend on the sugar in actual use.

- If none of the subsections below applies, then $match'(x, y)$ is true if and only if x and y have identical functors, identical arities, and $match(x_i, y_i)$ is true recursively for all their corresponding children x_i, y_i . Note that $match'$ recurses to $match$, not itself, so there is a chance for head matching on the children even if it was not relevant to the parent nodes.

Match anything The value of $match'(?), y)$ is always true. Thus, `?` can be used as a wildcard in `idsgrep` patterns to match an entire subtree regardless of its structure. Mnemonic: question mark is a shell wildcard for matching a single character. The verbose ASCII name for “`(?)`” is “(anything).”

Match anywhere The value of $match'(...x, y)$ is true if and only if there exists any subtree of y (including the entirety of y) for which $match(x, y)$ is true. In other words, this will look for an instance of x anywhere inside y regardless of nesting level. Mnemonic: three dots suggest omitting a variable-length sequence, in this case the variable-length chain of ancestors above x . The verbose ASCII name for “`...`” is “`.anywhere..`”

Match children in any order The value of $match'(*.x, y)$ is true if and only if there exists a permutation of the children of y such that $match(x, y')$ is true of the resulting modified y' . For instance, `*[a]bc` matches both `[a]bc` and `[a]cb`. This is obviously a no-operation (matches simply if x matches y , as if the asterisk were not applied) for trees of arity less than two. Mnemonic: asterisk is a general wildcard, and this is a general matching operation. The verbose ASCII name for “`.*`” is “`.unord..`”

NOT The value of $match'(!.x, y)$ is true if and only if $match(x, y)$ is false. It matches any tree *not* matched by x alone. Mnemonic: prefix exclamation point is logical NOT in many programming languages. The verbose ASCII name for “`!.x`” is “`.not..`”

AND The value of $match'([&]xy, z)$ is true if and only if $match(x, z) \wedge match(y, z)$. In other words, it

matches all trees that are matched by both x and y ; the set of strings matched by $[&]xy$ is the intersection of the sets matched by x and by y . Mnemonic: ampersand is logical or bitwise AND in many programming languages. The verbose ASCII name for “[&]” is “[and].”

OR The value of $match'([|]xy, z)$ is true if and only if $match(x, z) \vee match(y, z)$. In other words, it matches all trees that are matched by at least one of x or y ; the set of strings matched by $[|]xy$ is the union of the sets matched by x and by y . Mnemonic: ASCII vertical bar is logical or bitwise OR in many programming languages. The verbose ASCII name for “[|]” is “[or].”

Literal tree matching If x and y both have heads, then the value of $match'(.=.x, y)$ is true if and only if those heads are identical. Otherwise, it is true if and only if x and y have identical functors, identical arity, and $match(x_i, y_i)$ is true for each of their corresponding children.

The effect of this operation is to ignore any special $match'()$ semantics of x 's functor; the trees are compared as if that functor were just an ordinary string, regardless of whether it might normally be special. Note that the full $match()$ is still done on the children with only the root taken literally; to do a completely literal match of the entire trees it is necessary to insert an additional copy of $.=.$ above every node in the matching pattern, or at least every node that would otherwise have a special meaning for $match'()$, and even then heads will continue to have their usual effect of overriding recursion.[‡] Mnemonic: equals sign suggests the literal equality that is being tested rather than the more complicated comparisons that might otherwise be used. The verbose ASCII name for $.=.$ is $.equal.$

For instance, this feature could allow searching for a unary tree whose functor actually is $!$, where just specifying such a tree directly as the matching pattern would instead (under the rule for “NOT” above) search for trees that do not match the only child of $!$. In the original application of searching kanji decomposition databases this operation is unlikely to be used because the special functors do not occur

[‡]It may be interesting to consider how one could write a pattern to test absolute identity of trees, with each node matching if and only if its head or lack thereof is identical to the desired target as well as the functors and arities matching and the same being true of all children.

anyway, but it seems important for potential applications of IDSgrep to more general tree-querying, because otherwise some reasonable things people might want to look for could not be found at all.

Associative matching The value of $match'(.@.x, y)$ is calculated as follows. Create a new EIDS tree x' , initially equal to x , which has the property that its root may be of unlimited arity. Then for every child of x' whose functor and arity are identical to the functor and arity of x , replace that child in x' with its children, in order. Repeat that operation until no more children of x' have functor and arity identical to the functor and arity of x . Compute y' from y by the same process. Then $match'(.@.x, y) = match(=.x', y')$.

This matching operator is intended for the case of three or more things combined using a binary operator that has, or can be said to sometimes have, an associative law. For instance, the kanji 息 could be described by “ $\square\square\angle\square\heartsuit$ ” ($\angle\square$ over \heartsuit) or by “ $\square\angle\square\square\heartsuit$ ” (\angle over $\square\square\heartsuit$). Unicode might encourage use of the ternary operator \boxdot for this particular case instead, but that does not cover all reasonably-occurring cases, and the default databases seldom if ever use the Unicode ternary operators.

The difference between the representations is sometimes useful information that the database *should* retain; for instance, in the case of Tsukurimashou, “ $\square\square\angle\square\heartsuit$,” “ $\square\angle\square\square\heartsuit$,” and “ $\boxdot\square\square\heartsuit$ ” would correspond to three very different stanzas of MetaPost source code, and the user might want a query that separates them. On the other hand, the user might instead have a more general query along the lines of “find three things stacked vertically with \heartsuit at the bottom” and intend that that should match both cases of binary decomposition. The at-sign matching operation is meant for queries that don't care about the order of binary operators; without it, matching will by default follow the tree structure strictly.

Note that even with $.@.$, IDSgrep will not consider binary operators in any way interchangeable with ternary ones; users must still use $.|.$ to achieve such an effect if desired. Although the at-sign is fully defined for all arities, it is only intended for use with binary trees. Note also that $.@.$ and $.*$ behave according to their definitions. Incautious attempts to use them together will often fail to have the desired effects, because the definitions do not include special exceptions that some users might intuitively

expect for these two operators happening to occur near each other. In a pattern like “`*@[a][a]bcd`,” `.*` will recognize `.@` as the functor of a unary tree and expand the single permutation of its one child, and so that pattern will match the same things as if the asterisk had not been present, namely “`[a][a]bcd`” and “`[a]b[a]cd`” but not, for instance, “`[a][a]dcb`.” In a pattern like “`@[a]b*[a]cd`,” `.@` will recognize `.*` as a different arity and functor from `[a]` and choose not to expand it in x' , with the result that that pattern matches the same things as if the at-sign had not been present, namely “`[a]b[a]cd`” and “`[a]b[a]dc`” but not “`[a][a]bcd`” nor “`[a][a]bdc`.”

When considered as an operation on trees, what `.@` does is fundamentally the same thing as the algebraic operation that considers $(a + b) + c$ equivalent to $a + (b + c)$, and for that reason it is called “associative” matching. The mnemonic for at-sign is that it is a fancy “a” for “associative.” The verbose ASCII name for “`.@`” is “`.assoc.`”

Regular expression matching If x and y both have heads, then $match'(.l.x, y)$ is true if and only if the head of x , considered as a regular expression, matches the head of y . If x and y do not both have heads, then $match'(.l.x, y)$ is true if and only if x and y have the same arity, the functor of x considered as a regular expression matches the functor of y , and $match(x_i, y_i)$ is true for each of their corresponding children. This operation is basically the same as the default matching operation, except that regular expression matching is used instead of strict equality for testing the heads and functors. Mnemonic: slash means regular expression matching in Perl. Verbose ASCII name: “`.regex.`”

Regular expression matching for the purposes of this operator is as defined by the Perl Compatible Regular Expressions library, in whichever version was linked with the `idsgrep` utility. Strings are passed into PCRE as UTF-8, and are guaranteed (because the EIDS parser decodes and re-encodes `idsgrep`’s input for internal use) to be valid UTF-8 when PCRE sees them regardless of user input; as such, PCRE is given the option flags that make it read UTF-8 without doing its own validity check. Use of the PCRE “`\C`” syntax for matching individual octets within UTF-8 is strongly not recommended. All other PCRE options are left to the defaults chosen when PCRE was compiled, even if those are silly. The character tables are PCRE’s “C locale” defaults, not generated at runtime from the current locale. Things like case

sensitivity can be controlled within the pattern using PCRE’s syntax for doing so. In the event that `idsgrep` was compiled without the PCRE library (which is not recommended, but is possible), or that PCRE was compiled without UTF-8 support, then an attempt to evaluate the slash operator will trigger a fatal error.

A matching pattern given to PCRE will have already passed through the EIDS parser, which removes one level of backslash escaping. The pattern may also have been passed as a command-line argument to `idsgrep` by a shell, which may have undone another level of backslash escaping. Thus, it may be necessary to escape characters as many as three times in order to match them literally with the slash operator. Each of these levels may differ from the others in terms of the escape sequences it supports and their exact meanings. In many cases it doesn’t really matter which level of processing evaluates the escaping. For instance, “`idsgrep "/(\t)"`,” (shell evaluates “`\t`,” EIDS and PCRE see a literal tab); “`idsgrep "/"(\\t)"`,” (shell removes one backslash, EIDS evaluates “`\t`,” PCRE sees a literal tab); and “`idsgrep "/"(\\\\t)"`,” (shell removes two backslashes, EIDS removes one, PCRE evaluates “`\t`”) will all match the same things. If it matters, however, then caution is necessary.

PCRE because of the limitations of its API effectively forbids zero bytes (U+0000) in its matching patterns, whereas EIDS allows them to exist within strings in general. The complexities of PCRE pattern syntax make it impractical for `idsgrep` to automatically escape zero bytes before passing the strings to PCRE; there are too many different cases possible for the context in which a zero byte might occur. Since the `idsgrep` utility takes its matching patterns from the Unix command line anyway, and Unix itself forbids literal zero bytes in command-line arguments, the case of literal zero bytes in a matching pattern can only occur when they are created deliberately by escape sequences at the level of the EIDS parser; and the simplest advice to users is “don’t do that!”

Python, which like EIDS allows strings to contain zero bytes but has PCRE bindings and so faces the same issue, briefly attempted to work around this PCRE API limitation by auto-escaping. They eventually gave it up as too complicated and confusing. The consequence of PCRE’s API design is that if the string given as a matching pattern contains a literal zero byte then the regular expression to be matched will consist of the prefix of the string up to but not

including the first zero byte; anything after that will be ignored. Zero bytes are, nonetheless, permitted in the matching subject, and PCRE can search for them, but not by means of literal zero bytes in the pattern. For instance, the PCRE syntax “\000” (or just “\0” if the next character will not be an octal digit) matches a zero byte. As discussed above, additional escaping might be needed to ensure that PCRE, and not EIDS nor the shell, interprets the backslash escape.

User-defined matching predicates It is assumed that by some out-of-band means, we have defined a family of functions $U_i()$ for i from 1 up to some k . These functions take EIDS trees as input and return Boolean values (hence “predicates”).

Then the value of $match'(\#.x, y)$ is determined as follows. First, an integer i is computed. If x has a head, its initial characters will be parsed as an ASCII decimal number using the C library’s `atoi(3)` function; i is the resulting value, if it is positive. If x has no head, the head of x cannot be parsed, or the head of x is parsed as zero or negative, then i is defined to be 1. Having defined i , if $U_i()$ exists then $match' = U_i(y)$. If $U_i()$ does not exist then $match'$ is false. Mnemonic: hash-mark is used for parameter substitution in languages such as TeX, and this matching operation causes the matching pattern to take something external (the user-defined predicate) as a parameter.

In the current version, the functions $U_i()$ are always defined using the “-f” command-line option (or its long-named equivalent) and correspond to the character coverage of TrueType or OpenType fonts. The predicate returns true if and only if y has a head consisting of a single Unicode character that is covered by the font.

Cooked output

The default mode of operation for the `idsgrep` command-line utility is that whenever a matching tree is detected, the exact sequence of bytes that were parsed to generate that tree (including no skipped whitespace before it, and all skipped whitespace after it but before the next tree) will be copied through to the output. This mode of operation is called “raw.” Raw mode is easy to understand, efficient, preserves distinctions like different kinds of brackets in the input, and is as analogous as reasonably possible to the operation of `grep`. However, preserving the exact input bytes may preserve invalid UTF-8, valid but weird EIDS syntax, or non-ASCII charac-

ters users may find difficult to type or display, that may have existed in the input. The “-c” (“--cooking”) command-line option provides a wide range of ways for `idsgrep` to generate new EIDS syntax of its own, guaranteed to be valid, from the internal representation generated by the parser. The cooked output modes force the output into a well-behaved format independent of what the input looked like. Input canonicalization (such as the translation from “[lr]” to “[]”) can also be controlled through this interface.

The “-c” option can be given a (lowercase ASCII Latin, unabbreviated) keyword as its argument, to select a preset output mode. That is the only recommended way to use this option. The available preset modes are as follows:

raw Raw mode: write out the exact input byte sequence that was parsed to generate the matching tree, *even if it is not valid UTF-8*. This is the default.

rawnc Raw with no canonicalization: raw mode output, but without the canonicalization transformation during input parsing.

ascii ASCII-only: all non-ASCII characters and ASCII control characters are replaced by escape sequences or subjected to the reverse of the input canonicalization transformation, to produce a result that should pass through most limited-character-set channels. Note that the plainest ASCII space (U+0020) is not escaped in this mode when EIDS syntax does not require it to be. This mode generally uses a lot of hexadecimal escapes and, in a dictionary-lookup context, may be useful for finding the hexadecimal code point value of an unknown character.

cooked Generic cooked mode: render trees as reasonably clean and appealing Unicode text similar but not necessarily identical to what appears in the pregenerated dictionary files. This will escape characters outside the Basic Multicharacter Plane; characters in all Private Use Areas; and any other characters that EIDS syntax *requires* must be escaped; but no others. It will choose an appropriate escaping method depending on the type of character. Generally, it will use black lenticular brackets for top-level heads, ASCII brackets elsewhere, and syntactic sugar and syrup to avoid brackets where possible (except for top-level heads).

indent Write trees on multiple lines with two-space indentation to show their structure as clearly as possible. One blank line (two newlines) between trees. In other ways this is similar to “cooked.”

If not given a preset keyword, “-c” can be given a string of ASCII decimal digits. The decimal-string interface allows precise control of how output syntax will be generated, but it is somewhat experimental, very complicated, and may change incompatibly in future versions of this software. Use of this feature is not recommended. Nonetheless, the remainder of this section will attempt to document it.

The format specifier may be up to twelve digits long. If it is shorter than that, it is taken as a prefix with unspecified digits copied from the default specifier, which is “100000013250” and equivalent to the “cooked” preset. The two raw presets are handled as special cases; of the remaining cooked presets, “ascii” is equivalent to “000000013551” and “indent” is equivalent to “100000223250.”

The first digit specifies the type of brackets to be used for the head of the root of the tree: 0 for “<>,” 1 for “**[]**,” or 2 for “**[]**.” The second digit specifies the type of brackets for the head of any non-root node, using the same code.

The third digit specifies the type of brackets for nullary functors: 0 for “(),” 1 for “(),” or 2 for “()” . Similarly, the fourth digit specifies the brackets for unary functors: 0 for “.”, 1 for “:,” or 2 for “.”; the fifth digit specifies the brackets for binary functors: 0 for “[],” 1 for “[],” or 2 for “[]”; and the sixth digit specifies the brackets for ternary functors: 0 for “{ },” 1 for “[],” or 2 for “[]” .

The seventh digit describes how to insert newlines and indentation to pretty-print the tree structure. If it is 0, that will not be done. If it is 8, trees will be pretty-printed using one tab character per level; the number eight is a mnemonic for the fact that people generally expect those to be equivalent to eight spaces each. Any other decimal digit specifies that many spaces per level.

The eighth digit specifies the separator printed between trees: 0 for a null byte (U+0000), 1 for a newline, 2 for two newlines, or 3 for no separator at all.

The ninth digit specifies the circumstances under which the sugary and syrupy features of EIDS syntax should be used. It is a sum of binary flags: add 4 to use a syrupy semicolon when possible at the top level; 2 to use a syrupy semicolon when possible at other

levels; and 1 to use sugary implicit brackets wherever possible.

The tenth digit specifies which characters should be escaped. Literal backslashes, and (within a bracketed string) literal instances of the close-bracket character that would otherwise end the string, must always be escaped. When the tenth digit is 0, those are the only characters that will be escaped. Other values add escaping for the following categories of characters, and do so cumulatively with each digit also escaping everything that would be escaped by all lesser digits.

- 1 Escape characters from the astral planes; that is, characters with code points greater than U+FFFF and thus outside the Basic Multilingual Plane.
- 2 Escape characters from the BMP Private Use Areas, U+E000 to U+F8FF. The other Private Use Areas are already escaped at level 1 by virtue of being outside the BMP.
- 3 Escape all non-ASCII characters (U+0080 and up) except the core Unified Han range (U+4E00 to U+9FFF).
- 4 Escape the core Unified Han range.
- 5 Escape the ASCII control characters (U+0000 to U+001F).
- 6 Escape closing brackets at the start of bracketed strings, which otherwise escape escaping because of a special case in the syntax definition.
- 7 Escape all characters. Depending on the value of the next digit, however, the ASCII Latin alphabet still might not be escaped.

The eleventh digit specifies *how* to escape whatever characters were selected for escaping by the tenth digit. The available values are as follows.

- 0 Use a single backslash followed by the literal character, only. The ASCII Latin alphabet cannot be escaped in this way and under this option, or options 1 or 5 which fall through to this case, will not be escaped at all. Since the literal characters remain in the text, this option is not suitable for sending output through any channel that is not clean for the full range of UTF-8 characters. However, unlike raw mode, this and all other cooked modes do guarantee to produce valid UTF-8, not arbitrary byte sequences.

- 1 Use a backslash-letter sequence for ASCII control characters U+0001 to U+001B, and otherwise follow option 0.
- 2 Use variable-length hexadecimal “\x{” sequences for all characters that are selected to escape. This syntax can escape any character.
- 3 Use two-digit “\xHH” sequences wherever possible (that is, for ASCII and ISO-8859-1 characters), four-digit “\XHHHH” sequences for other characters on the Basic Multilingual Plane, and variable-length hexadecimal sequences otherwise.
- 4 Use four-digit “\XHHHH” sequences wherever possible (that is, for all characters on the BMP), and variable-length hexadecimal sequences otherwise.
- 5 Attempt to choose the simplest type of escape for each character depending on its value, just like option 3 except with backslash-letter escapes where possible (U+0001 to U+001B) and backslash-literal escapes for ASCII non-control characters (U+0020 to U+007E excluding the Latin alphabet). The ASCII Latin alphabet will not be escaped at all under this option.

The twelfth digit specifies canonicalization processing; that is, the translations on both input and output between alphabetic functor aliases like “(anything)” and their symbolic equivalents like “(?)”. Note that in all cases the symbolic versions are the matching operators; if you disable input canonicalization and enter a matching pattern of “(anything)” it will be matched as an ordinary nullary functor containing a string of eight ASCII letters, not as the match-anything operator which is always named “(?)”. The digit value is a sum of binary flags: add 4 to *disable* the default transformation of alphabetic aliases to symbolic names on input; plus 2 to enable a translation from alphabetic aliases to symbolic names on output, which is generally only meaningful if 4 was selected; plus 1 to enable a transformation from symbolic names back to alphabetic aliases on output.

Bibliography

- [1] CHISE project. Online <http://www.chise.org/>.
- [2] GlyphWiki. Online <http://en.glyphwiki.org/wiki/GlyphWiki:MainPage>.
- [3] Ulrich Apel. KanjiVG. Online <http://kanjivg.tagaini.net/>.
- [4] Jim Breen. The EDICT dictionary file. Online <http://www.csse.monash.edu.au/~jwb/edict.html>.
- [5] Jim Breen. WWWJDIC: Online Japanese Dictionary Service. Online <http://www.csse.monash.edu.au/~jwb/cgi-bin/wwwjdic.cgi>.
- [6] Alexandre Duret-Lutz. Using GNU Auto-tools. Online <http://www.lrde.epita.fr/~adl/dl/autotools.pdf>.
- [7] Free Software Foundation. GNU Grep 2.9. Online <http://www.gnu.org/software/grep/manual/grep.html>.
- [8] Jason Katz-Brown. The Kiten Handbook, revision 1.2. Online <http://docs.kde.org/development/en/kdeedu/kiten/index.html>.
- [9] 守岡知彦 [Morioka Tomohiko]. UTF-2000 プロジェクト [The UTF-2000 Project]. 漢字と情報 [Kanji and Information], (2):4–6, March 2001. In Japanese. Online <http://www.kanji.zinbun.kyoto-u.ac.jp/publications/kanji-and-info-2.pdf>.
- [10] Jon Postel. Transmission Control Protocol. RFC 793 (Standard), September 1981. Online <http://www.ietf.org/rfc/rfc793.txt>.
- [11] Julian Seward and Nicholas Nethercote. Using Valgrind to detect undefined value errors with bit-precision. In *USENIX Annual Technical Conference, General Track*, pages 17–30. USENIX, 2005.
- [12] Matthew Skala. Tsukurimashou Font Family and IDSgrep. Online <http://tsukurimashou.sourceforge.jp/>.
- [13] Matthew Skala. Tsukurimashou Github repository. Online <http://github.com/mskala/Tsukurimashou>.
- [14] Unicode Consortium. Ideographic description characters. In *The Unicode Standard, Version 6.0.0*, section 12.2. The Unicode Consortium, Mountain View, USA, 2011. Online <http://www.unicode.org/versions/Unicode6.0.0/ch12.pdf>.
- [15] Ben Wing et al. XEmacs: The next generation of Emacs. Online <http://www.xemacs.org/>.