

# System-Programmierung

## 9: Sockets

CC BY-SA, Thomas Amberg, FHNW  
(soweit nicht anders vermerkt)

# Ablauf heute

$\frac{2}{3}$  Vorlesung,

$\frac{1}{3}$  Hands-on,

Feedback.

Slides, Code & Hands-on: [tmb.gr/syspr-9](https://tmb.gr/syspr-9)



# Sockets

*Sockets* sind ein IPC Mechanismus um zwischen zwei Programmen Daten auszutauschen, die beide auf dem selben Host oder durch ein Netzwerk verbunden sind.

Die erste Implementierung des *Socket API* erschien 1983 mit 4.2BSD, deshalb auch "Berkeley Sockets".

Heute wird diese Schnittstelle für UNIX und Internet Sockets auf +/- allen Betriebssystemen unterstützt.

# Socket Verwendung

In einem typischen *Client-Server* Szenario nutzen Programme bzw. Anwendungen Sockets wie folgt:

Beide, Client und Server, kreieren einen Socket.

Der Server bindet seinen Socket auf eine wohlbekannte Adresse, so dass der Client ihn findet.

Kommunikation erfolgt uni- oder bidirektional.

# Socket Domänen

Die *Domäne* (communication domain) eines Sockets bestimmt, wie eine Socket Adresse aussieht, und ob lokal oder über ein Netzwerk kommuniziert wird.

Heutige Betriebssysteme unterstützen mindestens die UNIX (*AF\_UNIX* bzw. *AF\_LOCAL*) Domäne auf dem Host, sowie die Domänen IPv4 (*AF\_INET*) und IPv6 (*AF\_INET6*) für *Internet Protocol* (IP) Netzwerke.

# Stream Sockets

*Stream Sockets (SOCK\_STREAM)* sind zuverlässige, bidirektionale, verbindungsorientierte Byte Streams.

*Zuverlässig:* Bytes kommen entweder genau so an wie gesendet, oder Sender erhält eine Fehler-Notifikation.

*Bidirektional:* Datenübertragung in beide Richtungen, wie zwei Pipes, aber über ein Netzwerk. Deshalb auch *verbindungsorientiert:* verbunden mit einem *Peer*.

# Datagram Sockets

*Datagram Sockets (SOCK\_DGRAM)* sind Message-basiert, verbindungslos und unzuverlässig.

*Verbindungslos* bedeutet, dass einzelne Messages verschickt werden, ohne dass eine Verbindung da ist.

*Unzuverlässig* heisst, Übertragung und Reihenfolge sind nicht garantiert, Mehrfachübertragung möglich.

# Socket System Calls\*

Der *socket()* System Call kreiert einen neuen Socket.

Mit *bind()* binden Server ein Socket an eine Adresse.

Mit *listen()* hört ein Server auf neue Verbindungen.

Mit *accept()* wird eine Verbindung angenommen.

Der *connect()* System Call erstellt eine Verbindung mit einem anderen Socket. (\*Linux: Library Calls.)



# Socket kreieren mit *socket()*

Socket kreieren mit Domäne *domain* und Typ *type*:

```
int socket( // liefert einen File Deskriptor
    int domain, // AF_UNIX oder AF_INET, AF_INET6
    int type, // SOCK_STREAM oder SOCK_DGRAM
    int protocol); // immer 0 für diese Typen
```

Im Fehlerfall liefert *socket()* *-1* und setzt *errno*.

# Socket an Adresse binden mit *bind()*

Socket *sock\_fd* an die Adresse *sock\_addr* binden:

```
int bind( // 0 bei Erfolg, sonst -1 und errno  
         int sock_fd, // von socket() erstellt  
         const struct sockaddr *sock_addr,  
         socklen_t sock_addr_len);
```

Die Adresse hat je nach Domain einen anderen Typ,  
UNIX Domain Sockets verwenden einen Pfadnamen,  
Internet Sockets eine IP Adresse und einen Port.

# Socket Adressen

Der Struct *sockaddr* ist ein generischer Platzhalter:

```
struct sockaddr {  
    sa_family_t sa_family; // AF_ Konstante  
    char sa_data[14]; // Länge variiert  
}
```

Der *sa\_family* Wert genügt, um *sa\_data* zu parsen.

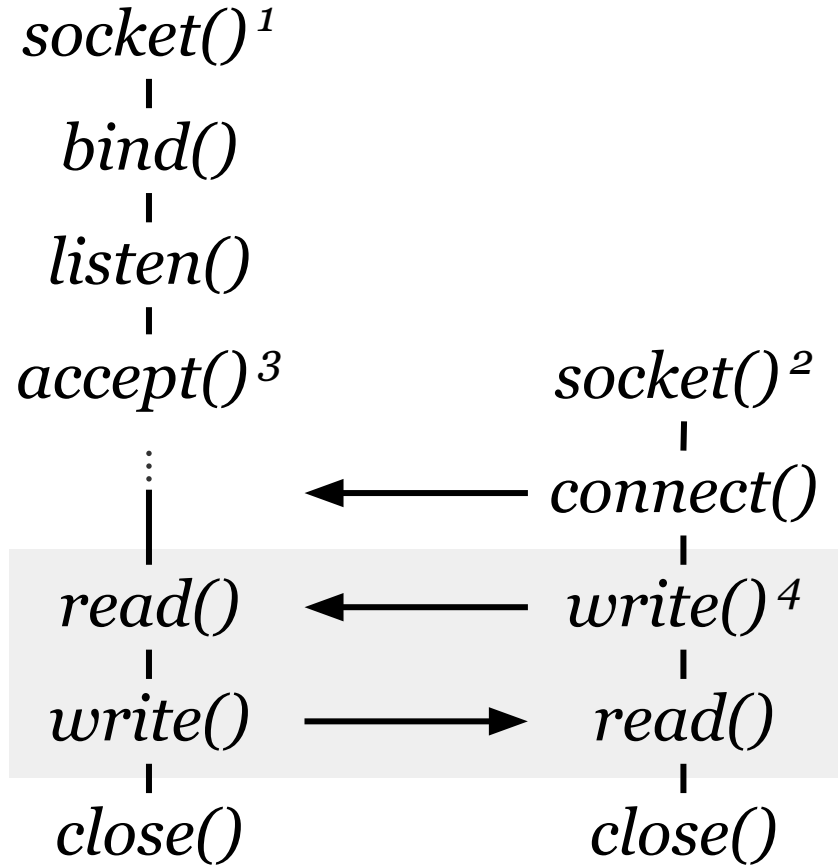
Der *sockaddr* Typ wird nur für Type-casts benutzt.

# Stream Sockets "Telefon" Analogie

*socket()* "installiert einen Telefonanschluss",  
*bind()* "löst eine Nummer", macht adressierbar,  
*listen()* "schaltet das Telefon ein", macht anrufbar.

*connect()* "ruft eine Nummer an", bzw. Adresse,  
*accept()* "nimmt einen eingehenden Anruf an",  
*send()* & *recv()* ist "reden & zuhören", bidirektional,  
*close()* "beide hängen auf am Ende des Anrufs".

# Stream Sockets Ablauf



<sup>1</sup> Server, passiver Socket.

<sup>2</sup> Client, aktiver Socket.

<sup>3</sup> Accept blockiert, bis zum Connect.

<sup>4</sup> Write kann von beiden Seiten initiiert werden, auch mehrfach.

# Auf Connections hören mit *listen()*

Auf eingehende Connections hören mit *listen()*:

```
int listen(int sock_fd, int backlog);
```

Muss vor *accept()* und *connect()* aufgerufen werden.

Der *backlog* Parameter bestimmt die Anzahl *pending* Connections, die von *accept()* angenommen werden.

Im Fehlerfall liefert *listen()* *-1* und setzt *errno*.

# Connections annehmen mit *accept()*

Eingehende Connections annehmen mit *accept()*:

```
int accept( // remote Socket fd, od. -1, errno  
    int sock_fd, // lokaler Socket File Deskr.  
    struct sockaddr *addr, // remote Adresse  
    socklen_t *addr_len); // Struct Grösse
```

Kreiert einen neuen Socket, der mit dem remote Peer / Client verbunden ist, der *connect()* aufgerufen hat.

Der Server Socket *sock\_fd* wird weiter verwendet. **n|w**

# Socket verbinden mit *connect()*

*connect()* verbindet zu einem Server bzw. Peer Socket:

```
int connect( // remote Socket fd, od. -1, errno  
            int sock_fd, // lokaler Socket File Deskript.  
            const struct sockaddr *addr, // remote Adr.  
            socklen_t addr_len); // Struct Grösse
```

Falls *connect()* einen Fehler liefert, Socket schliessen mit *close()* und neuen Socket kreieren mit *socket()*.



# Lesen und Schreiben mit *read()/write()*

Sockets sind bidirektional, beide Seiten können mit *read()/write()* oder *send()/recv()* lesen/schreiben.

Das Verhalten ist vergleichbar mit dem von Pipes, falls ein Ende geschlossen wird, kommt am anderen *EOF* raus bei *read()*, bzw. *EPIPE* bei *write()*, wenn zuvor das *SIGPIPE* Signal ignoriert worden ist.

Mit *close()* schliesst man eine Verbindung.

# Datagram Socket "Paketpost" Analogie

*socket()* "installiert einen Briefkasten",

*bind()* "weist dem Briefkasten eine Adresse zu".

*sendto()* "schickt ein Paket an einen Empfänger",

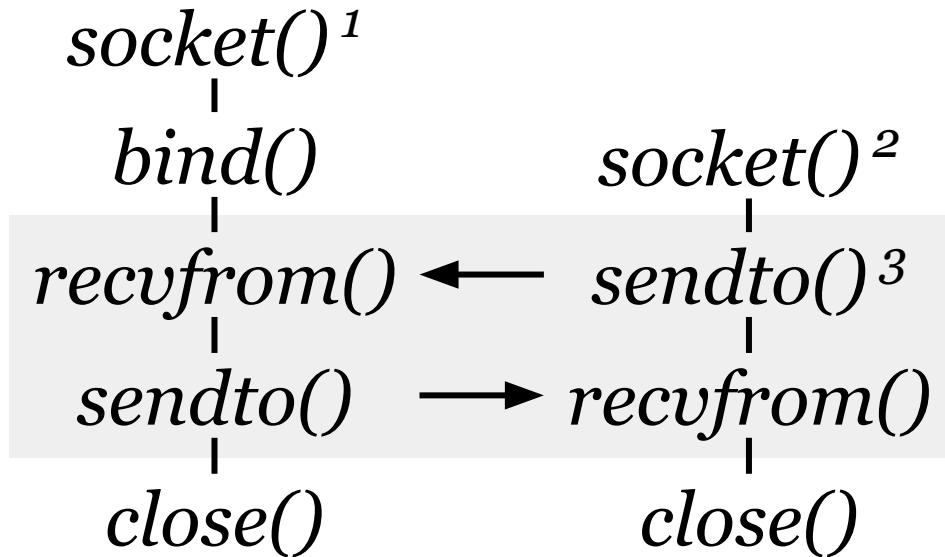
*recvfrom()* "wartet auf ein Paket, sieht Absender".

Die Pakete kommen in beliebiger Reihenfolge an.

*close()* "gibt den Briefkasten wieder frei".

# Datagram Sockets Ablauf

Bei Datagram Sockets entfällt *listen()* und *accept()*, sowie *connect()*, da diese verbindungslos sind.



<sup>1</sup> Server, passiver Socket.

<sup>2</sup> Client, aktiver Socket.

<sup>3</sup> Auch mehrfach und in beide Richtungen, weil *recvfrom()* die Adresse des Absenders liefert. **n|w**

# Datagram empfangen mit *recvfrom()*

Datagram empfangen, blockierend, mit *recvfrom()*:

```
ssize_t recvfrom( // Resultat wie bei read()
    int socket_fd, // Socket FD wie bei read()
    void *restrict buf, // wie bei read()
    size_t buf_len, // wie bei read()
    int flags, // 0, oder Socket spezifisch
    struct sockaddr *restrict source_addr,
    socklen_t *restrict source_addr_len);
```

# Datagram senden mit *sendto()*

Datagram senden an *dest\_addr* mit *sendto()*:

```
ssize_t sendto( // Resultat wie bei write()
    int sock_fd, // Socket FD wie bei write()
    const void *buf, // wie bei write()
    size_t buf_len, // wie bei write(), 0 ist OK
    int flags, // 0, oder Socket spezifisch
    const struct sockaddr *dest_addr,
    socklen_t dest_addr_len);
```

# UNIX Domain Sockets

UNIX Domain Sockets erlauben die Kommunikation zwischen zwei Prozessen auf demselben Host System.

Es gibt es sowohl Stream als auch Datagram Sockets.

Der Zugriff darauf ist über File Permissions geregelt.

*socketpair()* kreiert ein UNIX Domain Socket Paar.

Linux bietet einen abstrakten Socket Namespace.

# UNIX Domain Socket Adressen

Struct für Socket Adresse in der UNIX Domain:

```
struct sockaddr_un {  
    sa_family_t sun_family; // Immer AF_UNIX  
    char sun_path[108]; // Null-terminierter  
                          // Socket File-Pfad  
};
```

Die max. Länge von *sun\_path* ist Plattform-abhängig.

Deshalb beim Zuweisen *strncpy()* verwenden.

# UNIX Domain Socket binden mit *bind()*

Socket *sock\_fd* an die Adresse *addr* binden:

```
struct sockaddr_un addr;  
memset(&addr, 0, sizeof(struct sockaddr_un));  
addr.sun_family = AF_UNIX;  
strncpy(addr.sun_path, "/tmp/mysock",  
        sizeof(addr.sun_path) - 1);  
int sock_fd = socket(AF_UNIX, SOCK_STREAM, 0);  
bind(sock_fd, (struct sockaddr *) &addr,  
      sizeof(struct sockaddr_un));
```



# UNIX Domain Socket *bind()* Details

Der File-Pfad *addr.sun\_path* muss schreibbar sein.

UNIX Domain Sockets sind im RAM, nicht auf Disk.

Bestehenden Pfad erneut binden gibt *EADDRINUSE*.

Sockets werden an genau einen File-Pfad gebunden.

File *open()* funktioniert nicht auf Socket File-Pfad.

Unbenutzte Sockets mit *remove()* entfernen.

# Hands-on, 15': UNIX Domain Sockets

Analysieren Sie diese Socket Beispiele bestehend aus:

Header `us_xfr.h`<sup>TLPI</sup>,      Header `ud_ucase.h`<sup>TLPI</sup>,  
Server `us_xfr_sv.c`<sup>TLPI</sup>,      Server `ud_ucase_sv.c`<sup>TLPI</sup>,  
Client `us_xfr_cl.c`<sup>TLPI</sup>.      Client `ud_ucase_cl.c`<sup>TLPI</sup>.

Builden Sie die Programme, und lassen Sie sie laufen.

Zeichnen Sie **Sequenzdiagramme** mit User, Client, Server, das den Ablauf / übertragene Daten zeigt.

# UNIX Domain Socket Permissions

File Permissions bestimmen, wer Lese- oder Schreib-Zugriff auf UNIX Domain Sockets bekommen kann.

*bind()* erzeugt einen Socket Eintrag im File-System, mit allen Permissions für *owner*, *group* und *other*.

Für *connect()* und *sendto()* ist Schreibzugriff nötig, zudem braucht es *execute* (Such-) Rechte in jedem Directory des Socket File-Pfads.

# UNIX Domain Datagram Sockets

*UNIX Domain Datagram Sockets* nutzen wie Stream Sockets File-Pfade als Adresse, z.B. */tmp/mysocket*.

UNIX Domain Datagram Sockets übertragen Daten-Pakete zuverlässig, sequentiell und ohne Duplikate, im Gegensatz zu *Internet Domain Datagram Sockets*.

Pakete die grösser sind, als der bei *recvfrom()* mitgegebene Buffer werden abgeschnitten empfangen.

# UNIX Domain Datagram Paket-Grösse

Die maximale Grösse hängt von der Konfiguration ab.

Der Wert kann über die Socket Option *SO\_SNDBUF* mit *set-* / *getsockopt()* gesetzt bzw. abgefragt werden.

Beim Setzen eines Werts wird dieser verdoppelt (!), um Platz für die interne "Buchhaltung" zu schaffen:

```
setsockopt(fd, ..., SO_SNDBUF, ..., n);  
getsockopt(fd, ..., SO_SNDBUF, ..., &m); // 2*n
```

# Socket Paar kreieren mit *socketpair()*

Unbenanntes Socket Paar kreieren mit *socketpair()*:

```
int socketpair( // 0 oder -1, errno
               int domain, // nur für UNIX Domain AF_UNIX
               int type, // SOCK_DGRAM oder SOCK_STREAM
               int protocol, // 0
               int sock_fd[2]); // zwei verbundene Sockets
```

Typischerweise gefolgt von *fork()*, wie bei *pipe()*.

Kein File-Pfad => "unsichtbar", bessere Security. 

# Linux Abstract Socket Namespace

Der *abstract* Namespace ist ein Linux-spezifisches Feature um UNIX Domain Sockets an einen Namen zu binden, der nicht im File-System kreiert wird.

Ohne Pfadname keine Kollisionen, kein *remove()*.

Um einen Namen im *abstract* Namespace zu kreieren, setzt man in *addr.sun\_path* den ersten *char* auf ' $\emptyset$ '.

(Kein gutes API Design, ging wohl nicht anders.)

# Internet Domain Sockets

Internet Domain *Stream Sockets* basieren auf dem *TCP* Protokoll. Sie bieten zuverlässige, bidirektionale Kommunikation mit Byte Stream Semantik.

Internet Domain *Datagram Sockets* basieren auf dem *UDP* Protokoll. Im Unterschied zu der UNIX Variante sind UDP Sockets nicht zuverlässig, garantieren keine Ordnung, es gibt Duplikate und "dropped packets".



# Netzwerk Byte Reihenfolge

Die *Network Byte Order* ist eine Konvention wie man Integer Werte in Bytes zerlegt und zwar "Big Endian".

Bei *Big Endian* schreibt man das *MSB* vor dem *LSB*:



Library Funktionen die IP Adressen ausgeben, liefern Resultate immer in Network Byte Order. Konstanten wie *INADDR\_ANY* müssen konvertiert werden.

# Byte Reihenfolge konvertieren

Konvertieren von Netzwerk zu Host Byte Order:

```
uint32_t ntohs(uint32_t netlong);  
uint16_t ntohs(uint16_t netshort);
```

Konvertieren von Host zu Netzwerk Byte Order:

```
uint32_t htonl(uint32_t hostlong);  
uint16_t htons(uint16_t hostshort);
```


Die Host Byte Reihenfolge kann je nach Hardware Plattform entweder Big oder Little Endian sein.

# Repräsentation von Daten

Nicht nur bei Adressen, auch bei allen anderen via ein Netzwerk gesendeten Daten ist das *Encoding* wichtig.

Bei TCP und UDP legt das die Anwendungsebene fest.

*HTTP* fordert z.B. *US-ASCII* für den Message Header, und via *Content-Type* beliebige Content Encodings.

Content vom Typ *application/json* würde z.B. gemäss JSON Standard mit UTF-8 Encoding übertragen. 

# IPv4 Internet Socket Adressen

IPv4 Internet Socket Adresse, z.B. 192.168.0.42

```
struct in_addr {  
    uint32_t s_addr; // Network Byte Order  
};  
  
struct sockaddr_in {  
    sa_family_t sin_family; // AF_INET  
    in_port_t sin_port; // Network Byte Order  
    struct in_addr sin_addr; // Internet Adresse  
};
```

# IPv6 Internet Socket Adressen

```
struct in6_addr {  
    unsigned char s6_addr[16]; // IPv6 address  
};
```

```
struct sockaddr_in6 {  
    sa_family_t sin6_family; // AF_INET6  
    in_port_t sin6_port; // Port Nummer  
    uint32_t sin6_flowinfo; // IPv6 Flow Info  
    struct in6_addr sin6_addr; // IPv6 Adresse  
    uint32_t sin6_scope_id; // Scope ID  
};
```

# Loopback und Wildcard Adressen

IPv4 Loopback 127.0.0.1 und Wildcard 0.0.0.0 Adr.:  
INADDR\_LOOPBACK, INADDR\_ANY

IPv6 Loopback (::1) und Wildcard (::) Adresse:  
in6addr\_loopback bzw. IN6ADDR\_LOOPBACK\_INIT,  
in6addr\_any bzw. IN6ADDR\_ANY\_INIT

# Internet Socket Adressen Konvertieren

Von Punkt-Notation zu Binärformat konvertieren:

```
int inet_pton( // Erfolg: 1, Fehler: 0 od. -1
    int addr_family, // AF_INET, AF_INET6
    const char *src, // IP Adr. in Punkt-Notation
    void *dst); // IP Adresse im Binärformat
```

Von Binärformat zu Punkt-Notation konvertieren:

```
const char *inet_ntop( // dst od. NULL, errno
    int addr_family, // AF_INET, AF_INET6
    const void *src, // IP Adresse im Binärformat
    char *dst, socklen_t size); // IP String n|w
```

# Host Lookup mit *getaddrinfo()*

Lookup von *host* und *service* (mit *hints*) liefert *result*:

```
int getaddrinfo( // 0 bei Erfolg, sonst != 0
    const char *host, // Hostname od. IP Adresse
    const char *service, // Name od. Port Nummer
    const struct addrinfo *hints, // Bsp. unten
    struct addrinfo **result); // Liste, ai_next
```

Nach Gebrauch, *addrinfo* Struct *result* freigeben:

```
void freeaddrinfo(struct addrinfo *result)
```



# Struct *addrinfo*

```
struct addrinfo { // hint* u. result, Rest = 0
    int ai_flags*; // Siehe Doku für AI_... Flags
    int ai_family*; // AF_UNSPEC, AF_INET(6)
    int ai_socktype*; // SOCK_STREAM, SOCK_DGRAM
    int ai_protocol*; // 0
    socklen_t ai_addrlen; // IP Adress-Länge
    struct sockaddr *ai_addr; // IP Adress-Struct
    char *ai_canonname; // Kanonischer Name
    struct addrinfo *ai_next; // "next" od. NULL
};
```

# Hands-on, 15': Internet Domain Sockets

Analysieren Sie dieses Socket Beispiel bestehend aus:

Header `i6d_ucose.h`<sup>TLPI</sup>,

Server `i6d_ucose_sv.c`<sup>TLPI</sup>,

Client `i6d_ucose_cl.c`<sup>TLPI</sup>.

Builden Sie die Programme, und lassen Sie sie laufen:

```
$ ./i6d_ucose_sv &
```

```
$ ./i6d_ucose_cl ::1 hello
```

# Hausaufgabe, 3h: Web Client und Server

Lesen Sie die Kapitel 4 bis 7 der HTTP Spezifikation  
<https://tools.ietf.org/html/rfc2616>

Lösen Sie die beiden folgenden Hands-on Aufgaben.

# Hands-on, 1h: Web Client `http_client.c`

Schreiben Sie einen Web Client *my\_http\_client.c*, der folgenden HTTP Request an den Host *tmb.gr*, Port 80 sendet, die Antwort liest, und auf *stdout* ausgibt:

```
"GET /syspr HTTP/1.1\r\n"
```

```
"Host: tmb.gr\r\n"
```

```
"\r\n"
```

Hinweis: HTTP nutzt TCP als Transport-Protokoll.

Länge der Antwort ist im *Content-Length* Header. 

# Hands-on, 1h: Web Server `http_server.!c`

Schreiben Sie einen Web Server *my\_http\_server.c*,  
der einkommende HTTP Requests auf Port 8080 liest  
und folgende Antwort zum Client / Browser sendet:

```
"HTTP/1.1 200 OK\r\n"
```

```
"Connection: close\r\n"
```

```
"Content-Length: 5\r\n"
```

```
"\r\n"
```

```
"hello"
```

# Feedback?

Gerne im [Slack](#) oder an [thomas.amberg@fhnw.ch](mailto:thomas.amberg@fhnw.ch)

Programmierfragen am besten schriftlich.

Sprechstunde auf Voranmeldung.

Slides, Code & Hands-on: [tmb.gr/syspr-9](https://tmb.gr/syspr-9)



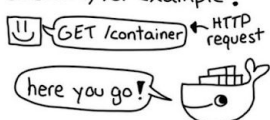
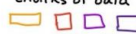





Julia Evans  
@b0rk

Following

## unix domain sockets

unix domain sockets JULIA EVANS  
@b0rk

<p>unix domain sockets are files.</p> <p><code>\$ file mysock.sock</code> socket</p> <p>the file's permissions determine who can send data to the socket</p>	<p>they let 2 programs on the <u>same computer</u> communicate.</p> <p>Docker uses Unix domain sockets, for example!</p> 	<p>There are 2 kinds of unix domain sockets:</p> <p><b>stream</b> like TCP! Lets you send a continuous stream of bytes</p> <p><b>datagram</b> like UDP! Send discrete chunks of data</p> 
<p><b>advantage 1</b></p> <p>Lets you use file permissions to restrict access to HTTP/ database services!</p> <p><code>chmod 600 secret.sock</code></p> <p>This is why Docker uses a unix domain socket</p> 	<p><b>advantage 2</b></p> <p>UDP sockets aren't always reliable (even on the same computer), unix domain datagram sockets <u>are</u> reliable! And won't reorder!</p> 	<p><b>advantage 3</b></p> <p>You can send a file descriptor over a unix domain socket. Useful when handling untrusted input files!</p> 

1:51 AM - 3 Apr 2018

Julia Evans  
@b0rk

programming and  
she/her. ❤️ @reci  
@bangbangcon

Montreal



Tweet

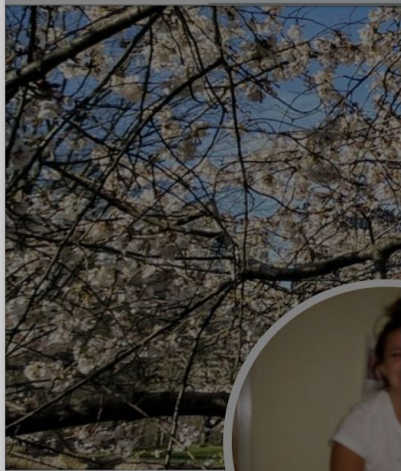


Following



Follow





jessie frazelle  
@jessfraz

A superhero with  
Keyser Söze of cc  
@Microsoft, Xoog  
maintainer. conta



jessie frazelle ✓ @jessfraz · Sep 22

You guys it took me 4 days to realize and find a fix for a bug so here's the deal: I needed to get a value from a bitfield so I calculated the offset from the closest neighbor... except the value it was returning did not make sense in the context... I'm like it's gotta be my math

4 2 55



jessie frazelle ✓

@jessfraz

Following

No it's not. It was the freaking kernel version of the struct I was looking up was different than the one I was using so the items were in a different order.

5:50 AM - 22 Sep 2018

54 Likes



8 54



Tweet your reply



jessie frazelle ✓ @jessfraz · Sep 22