



دانشگاه صنعتی شریف
دانشکده مهندسی کامپیوتر

گزارش پژوهش
درس سیستم عامل

پشتیبانی از شبکه در سطح هسته در سیستم عامل **XV**

نگارش

فاتیما تیمارچی - ۰۲۱۰۵۸۰۲

آیه صابری - ۰۲۱۰۶۱۴۵

متین باقری - ۰۲۱۰۵۷۲۷

حورا عابدین - ۰۱۱۰۶۲۰۹

استاد

دکتر حسین اسدی

۱۴۰۴ بهمن ماه

فهرست مطالب

۱	۱	۱	۱	مقدمه
۱	۱	۱	۱	۱-۱ تعریف مسئله
۲	۲	۲	۲	۲-۱ اهداف پروژه
۲	۲	۲	۲	۲-۲ ساختار گزارش
۳	۳	۳	۳	۳ گام اول: تکمیل بخش‌های اولیه درایور E1000
۳	۳	۳	۳	۳-۱ مقداردهی اولیه ring RX
۳	۳	۳	۳	۳-۲ مقداردهی اولیه ring TX
۴	۴	۴	۴	۴-۱ تکمیل تابع transmit_1000e
۴	۴	۴	۴	۴-۲ تکمیل تابع recv_1000e
۵	۵	۵	۵	۵ گام دوم: پیاده‌سازی مسیر ارسال (TX)
۶	۶	۶	۶	۶ گام سوم: پیاده‌سازی مسیر دریافت (RX)
۷	۷	۷	۷	۷ گام چهارم: کار با register های کارت شبکه
۸	۸	۸	۸	۸ گام پنجم: پیاده‌سازی دریافت UDP در net.c
۹	۹	۹	۹	۹ گام ششم: تکمیل پردازش UDP در ip_rx
۱۱	۱۱	۱۱	۱۱	۱۱ گام هفتم: مدیریت آزادسازی های buffer شبکه
۱۲	۱۲	۱۲	۱۲	۱۲ سناریوی تست شبکه

۱۰ بخش امتیازی: پشتیبانی از ICMP و گزارش خطای

۱۴

فهرست تصاویر

۱-۹ ارسال بسته‌های UDP از سمت host برای تست دریافت و محدودیت صفحه	۱۳
۲-۹ خروجی پس از دریافت بسته‌های UDP	۱۳

فصل ۱

مقدمه

در این پروژه، پشتیبانی شبکه را در سیستم عامل آموزشی تکمیل کردیم. تمرکز اصلی ما پیاده‌سازی درایور کارت شبکه‌ی E1000 در سطح kernel و اتصال آن به stack network سیستم عامل بود. این کارت شبکه در محیط QEMU شبیه‌سازی می‌شد، اما از دید kernel رفتاری مشابه یک کارت شبکه‌ی واقعی دارد. به همین دلیل لازم بود بخش‌های مربوط به TX و RX register را کامل کنیم، با register های کارت شبکه که در حافظه در دسترس هستند کار کنیم و interrupt های ایجاد شده را مدیریت کنیم تا ارتباط بین سخت‌افزار و هسته به درستی برقرار شود.

۱-۱ تعریف مسئله

مسئله‌ی اصلی پروژه ایجاد یک مسیر کامل برای ارسال و دریافت packet های شبکه داخل سیستم عامل بود. در بخش ارسال (TX)، زمانی که stack network قصد ارسال یک packet را داشت، آدرس buffer داده را داخل ring TX قرار دادیم و با به روزرسانی register مربوطه مانند TDT، کارت شبکه را از وجود داده‌ی جدید مطلع کردیم. پس از این‌که کارت شبکه ارسال را کامل کرد و بیت وضعیت (DD) را فعال کرد، buffer را با kfree آزاد کردیم تا مدیریت حافظه به درستی انجام شود.

در مسیر دریافت (RX)، کارت شبکه پس از دریافت packet از شبکه، داده را از طریق DMA داخل buffer مشخص شده می‌نوشت و یک interrupt ایجاد می‌کرد. در تابع مربوط به دریافت، ring RX را بررسی کردیم، های packet جدید را به لایه‌ی بالاتر تحویل دادیم و سپس برای هر slot مصرف شده یک buffer جدید با kalloc اختصاص دادیم تا حلقه برای دریافت packet های بعدی آماده بماند.

۲-۱ اهداف پروژه

هدف این پروژه تکمیل پشتیبانی UDP در سیستم عامل و فراهم کردن امکان دریافت داده در user space بود. برای این کار، syscall های bind ، recv و unbind را پیاده سازی کردیم تا هر process بتواند روی یک port مشخص منتظر دریافت داده بماند. در لایه IP بررسی کردیم که دریافتی از نوع UDP باشد و port مقصد آن قبلاً bind شده باشد. در صورت برقرار بودن این شرایط، packet در queue مربوط به همان port قرار می گیرد. برای هر port حداقل ۱۶ packet در queue نگه می داریم تا مصرف حافظه کنترل شود و اگر صفحه pر باشد، packet جدید drop می شود. همچنین در تمام مسیر دقت کردیم هر buffer که دیگر مورد نیاز نیست، چه در حالت تحويل به کاربر و چه در حالت drop شدن، آزاد شود.

۳-۱ ساختار گزارش

ساختار این گزارش مطابق با گام های تعریف شده در صورت پروژه تنظیم شده است. در هر بخش، توضیح پیاده سازی همان گام در کد ارائه شده و روند تکمیل پروژه از درایور کارت شبکه تا پشتیبانی UDP در فضای کاربر به صورت مرحله به مرحله بررسی شده است.

فصل ۲

گام اول: تکمیل بخش‌های اولیه درایور E1000

در گام اول، فایل kernel/e1000.c که بخش اصلی درایور کارت شبکه در آن قرار دارد را کامل کردیم تا زیرساخت ارسال و دریافت packet در سطح kernel آماده شود.

۱-۲ مقداردهی اولیه RX ring

در بخش دریافت، ابتدا RX ring را به صورت کامل مقداردهی کردیم تا کارت شبکه هنگام دریافت packet فضای معتبر برای نوشتمندانه داشته باشد. برای این کار:

- برای هر خانه از RX ring یک buffer با استفاده از () kalloc ساختیم و آدرس آن را در فیلد addr همان descriptor قرار دادیم تا کارت بتواند از طریق DMA داده را مستقیم داخل حافظه بنویسد.
- آدرس پایه‌ی RX و اندازه‌ی آن را در registerهای مربوطه تنظیم کردیم تا کارت شبکه بداند حلقه در چه بخشی از حافظه قرار دارد.
- مقادیر اولیه‌ی RDH و RDT را تنظیم کردیم تا وضعیت ابتدایی حلقه مشخص باشد و دریافت بتواند شروع شود.

۲-۲ مقداردهی اولیه TX ring

در بخش ارسال هم TX ring را آماده کردیم تا بتواند های packet خروجی را مدیریت کند. در این قسمت:

- آدرس پایه‌ی TX ring را در registerهای مربوطه قرار دادیم تا کارت شبکه بتواند descriptorهای ارسال را بخواند.

- طول حلقه را تنظیم کردیم و مقادیر اولیه‌ی TDT و TDH را مشخص کردیم.
- وضعیت اولیه‌ی descriptor های TX را تنظیم کردیم تا هنگام اولین ارسال، هر خانه در وضعیت مشخصی قرار داشته باشد.

۳-۲ تکمیل تابع transmit_1000e

در تابع transmit_1000e بخش‌هایی را اضافه کردیم تا وقتی network stack درخواست ارسال دارد، packet به درستی داخل TX ring قرار بگیرد. در این تابع:

- مقدار فعلی TDT را خواندیم تا اندیس descriptor بعدی مشخص شود.
- وضعیت آن descriptor را بررسی کردیم تا مطمئن شویم هنوز در حال استفاده نیست.
- آدرس buffer ارسالی را در فیلد addr قرار دادیم و مقدار length را تنظیم کردیم.
- بیت‌های کنترلی لازم برای ارسال (مانند EOP و RS) را فعال کردیم.
- در پایان، مقدار TDT را جلو بردیم تا کارت شبکه متوجه شود داده‌ی جدید برای ارسال آماده شده است.

۴-۲ تکمیل تابع recv_1000e

تابع recv_1000e را هم کامل کردیم تا packet های دریافتی از روی RX ring خوانده شوند و به لایه‌ی بالاتر تحویل داده شوند. منطق این تابع به این صورت پیاده‌سازی شد:

- حلقه‌ی RX ring را بررسی کردیم و descriptor هایی که وضعیت دریافت آنها فعال شده بود شناسایی کردیم.
- buffer مربوط به هر packet جدید را به rx_net() تحویل دادیم تا پردازش ادامه پیدا کند.
- بعد از تحویل، یک buffer جدید با kalloc() ساختیم و جایگزین قبلی کردیم تا حلقه برای دریافت‌های بعدی آماده بماند.
- در نهایت RDT را به روزرسانی کردیم تا کارت شبکه بتواند دوباره از آن خانه استفاده کند.

فصل ۳

گام دوم: پیاده‌سازی مسیر ارسال (TX)

در گام دوم، ارسال بسته‌های شبکه در درایور کارت شبکه را پیاده کردیم. زمانی که پشته‌ی شبکه در فایل net.c نیاز به ارسال یک packet دارد، با فراخوانی تابع `transmit_1000e`، پوینتر buffer داده به درایور منتقل می‌شود و درایور آن را داخل TX ring قرار داده و کارت شبکه را مطلع می‌کند.

در ابتدای تابع `transmit_1000e`، مقدار فعلی register مربوط به TDT خوانده می‌شود تا مشخص شود کدام خانه از TX ring برای ارسال بعدی استفاده خواهد شد. این ایندکس تعیین می‌کند داده‌ی جدید در کدام descriptor قرار می‌گیرد و مدیریت صحیح آن از descriptor overwrite شدن descriptor های در حال استفاده جلوگیری می‌کند. سپس وضعیت descriptor بررسی می‌شود و بیت E1000_TXD_STAT_DD چک می‌شود تا مطمئن شویم ارسال قبلی کامل شده و آن خانه آزاد است.

در ادامه، آدرس buffer ارسالی که در فایل net.c با `kalloc` تخصیص داده شده در فیلد `addr` قرار می‌گیرد و طول packet در فیلد `length` تنظیم می‌شود. سپس بیت‌های کنترلی descriptor مقداردهی می‌شوند؛ بیت EOP برای مشخص کردن پایان packet و بیت RS برای درخواست گزارش وضعیت فعال می‌شود. فعال بودن RS باعث می‌شود کارت شبکه پس از اتمام ارسال، بیت DD را تنظیم کند.

پس از تکمیل descriptor، مقدار TDT یک واحد جلو برد شده و در register نوشته می‌شود. با این کار، کارت شبکه متوجه آماده بودن یک descriptor جدید شده و packet را از طریق DMA ارسال می‌کند.

در نهایت، آزادسازی buffer تنها پس از فعال شدن بیت DD انجام می‌شود و با `kfree` حافظه آزاد می‌گردد. این ترتیب مانع آزاد شدن زودهنگام حافظه و بروز خطأ در ارسال می‌شود.

فصل ۴

گام سوم: پیاده‌سازی مسیر دریافت (RX)

در این گام، منطق دریافت بسته‌ها در درایور کارت شبکه را تکمیل کردیم. زمانی که کارت شبکه یک packet را دریافت می‌کند، آن را از طریق DMA در حافظه‌ای می‌نویسد که آدرس آن در فیلد addr یکی از descriptorهای RX ring قرار دارد. پس از نوشتن داده، کارت شبکه با فعال کردن بیت وضعیت descriptor و ایجاد interrupt، kernel را از ورود داده‌ی جدید مطلع می‌کند. پردازش این داده درتابع recv_1000e انجام می‌شود.

در ابتدای تابع recv_1000e، اندیس بررسی از روی مقدار register مربوط به RDT محاسبه می‌شود. سپس descriptorهای بعدی در RX ring بررسی می‌شوند و در هر مرحله بیت وضعیت چک می‌شود. اگر بیت وضعیت فعال نشده باشد، یعنی داده‌ی جدیدی وجود ندارد و پردازش متوقف می‌شود. در غیر این صورت، buffer مربوط به descriptor و طول داده از فیلد های addr و length خوانده می‌شود.

در این مرحله buffer دریافتی مستقیماً به تابع rx_net() در فایل net.c تحویل داده می‌شود تا پردازش در لایه‌های بالاتر (Ethernet و IP و سپس UDP) ادامه پیدا کند. این کار باعث می‌شود داده بدون کپی اضافی وارد network stack شود.

پس از تحویل packet، همان slot از RX ring باید دوباره برای دریافت آماده شود. برای این کار، یک buffer جدید با kalloc() تخصیص داده شده و آدرس آن جایگزین قبلی می‌شود. سپس بیت وضعیت descriptor پاک می‌شود و مقدار RDT به روزرسانی می‌شود تا کارت شبکه بتواند دوباره از آن خانه استفاده کند.

فصل ۵

گام چهارم: کار با register های کارت شبکه

در این گام، ارتباط مستقیم درایور با register های کنترلی کارت شبکه را کامل کردیم. تا این مرحله TX و RX ring در حافظه بدرسی کار می‌کردند، اما برای اینکه کارت شبکه متوجه تغییرات شود، لازم بود مقدار بعضی از register ها را به صورت دقیق بخوانیم و به روزرسانی کنیم. دسترسی به این register ها از طریق آرایه‌ی سراسری regs در فایل e1000.c در مسیر ارسال

در مسیر ارسال، بعد از اینکه در تابع transmit_1000e descriptor مقداردهی شدند، مقدار [E1000_TDT] reg جلو برد شد. این همان لحظه‌ای است که کارت شبکه متوجه می‌شود یک descriptor جدید آماده ارسال است. اگر این مقدار نوشته نشود، حتی با تنظیم descriptor کامل هم packet ارسال نخواهد شد.

در مسیر دریافت نیز داخل تابع recv_1000e بعد از اینکه packet پردازش شد و buffer جایگزین گردید، مقدار [E1000_RDT] reg به روزرسانی شد. این کار باعث می‌شود کارت شبکه بداند آن خانه از ring RX دوباره آماده استفاده است. اگر RDT تغییر نکند، دریافت در همان موقعیت متوقف می‌شود.

در کل، در بخش‌های زیر با register ها کار داریم:

- خواندن [E1000_TDT] reg برای تعیین محل قرار دادن descriptor جدید در TX ring
- نوشتن مقدار جدید در [E1000_TDT] reg بعد از آماده‌سازی descriptor در transmit_1000e
- استفاده از [E1000_RDT] reg برای مشخص کردن موقعیت بررسی در RX ring داخل recv_1000e
- به روزرسانی [E1000_RDT] reg بعد از جایگزینی buffer جدید در مسیر دریافت

فصل ۶

گام پنجم: پیاده‌سازی دریافت UDP در c.net.c

در این گام، دریافت بسته‌های UDP در فایل kernel/net.c را تکمیل کردیم تا فرایندهای کاربر بتوانند داده‌های دریافتی را از طریق syscall دریافت کنند. تا این مرحله، packet ها از کارت شبکه وارد() می‌شدند، اما مکانیزمی برای نگهداری و تحويل آنها به user space وجود نداشت.

اول ساختاری برای نگهداری وضعیت هر port اضافه شد. برای هر port، یک صف از packet های دریافتی در نظر گرفته شد تا داده‌ها تا زمان فراخوانی recv در kernel باقی بمانند. برای جلوگیری از مصرف بیش از حد حافظه، برای هر port حداقل ۱۶ packet در صف نگه داشته می‌شود.

در پیاده‌سازیتابع bind، زمانی که یک port مشخص را ثبت می‌کند، وضعیت آن port در ساختار داخلی net.c علامت‌گذاری می‌شود تا packet های با آن مقصد ذخیره شوند. اگر قبل از ثبت شده باشد، فراخوانی با خطأ مواجه می‌شود. در مقابل، در unbind وضعیت همان port آزاد می‌شود و دیگر packet جدیدی برای آن ذخیره نخواهد شد.

در ادامه، در مسیر دریافت داخل (ip_rx(), بعد از اینکه مشخص شد packet از نوع UDP است، شماره‌ی port مقصد از header استخراج می‌شود. سپس بررسی می‌شود که آیا این port قبلاً توسط bind ثبت شده است یا خیر. اگر ثبت نشده باشد، packet آزاد می‌شود. اگر ثبت شده باشد و صف آن port هنوز ظرفیت داشته باشد، packet در صف قرار می‌گیرد؛ در غیر این صورت drop packet می‌شود تا از پر شدن حافظه جلوگیری شود.

در پیاده‌سازی recv، زمانی که process فراخوانی انجام می‌دهد، اگر در صف مربوط به port داده‌ای موجود باشد، اولین packet از صف برداشته شده و محتوای آن به فضای کاربر کپی می‌شود. پس از تحويل داده، buffer مربوطه با kfree آزاد می‌شود. اگر صف خالی باشد، رفتار مطابق طراحی پروژه انجام می‌شود (مانند بازگشت مقدار مناسب یا انتظار).

فصل ۷

گام ششم: تکمیل پردازش UDP در ip_rx

تمام تغییرات این گام در تابع ip_rx در فایل kernel/net.c انجام شد. این تابع برای هر بسته IP که از net_rx() وارد سیستم می‌شود فرآخوانی می‌شود.

در ابتدای ip_rx، ابتدا طول بسته بررسی می‌شود که حداقل شامل header اترنت و UDP است. سپس طول واقعی IP header از روی فیلد ip_vhl برداشته می‌شود تا بتوانیم offset دقیق header را محاسبه کنیم. اگر طول بسته کمتر از مقدار لازم باشد، buffer همانجا با kfree آزاد می‌شود.

بعد از تأیید طول، نوع پروتکل بررسی می‌شود و تنها در صورتی که ip_p == IPPROTO_UDP باشد، پردازش ادامه پیدا می‌کند. سپس pointer UDP header مربوط به با استفاده از offset محاسبه شده ساخته می‌شود. برای خواندن فیلدهای چندبایتی مانند شماره پورت و طول، از توابع ntohs() و ntohl() استفاده شده است تا تبدیل از ترتیب بایت شبکه به ترتیب بایت پردازنده به درستی انجام شود.

در این بخش، شماره پورت مقصید استخراج شده و اعتبار آن بررسی می‌شود. اگر پورت خارج از محدوده تعریف شده باشد یا قبلاً توسط bind() ثبت نشده باشد، بسته کنار گذاشته می‌شود. همچنین اگر صفت مربوط به آن پورت به حداکثر ظرفیت خود رسیده باشد، بسته drop می‌شود تا از مصرف بیش از حد حافظه جلوگیری شود.

سپس کنترل‌های زیر در ip_rx انجام می‌شود:

- بررسی حداقل طول بسته برای دسترسی امن به IP و UDP header
- بررسی نوع پروتکل و اطمینان از UDP بودن بسته
- تبدیل ترتیب بایت برای sport، dport و آدرس مبدأ
- کنترل معابر بودن پورت و ظرفیت صفت آن

در صورت معابر بودن شرایط، بسته در صفت مربوط به آن پورت ذخیره می‌شود. (ساختار صفت در

ابتدای فایل net.c تعریف شده است). داخل ip_rx_lock پس از گرفتن lock مربوط به آن پورت، محل ذخیره با استفاده از ایندکس UDPQSIZE % w تعیین می شود. سپس اطلاعات بسته شامل آدرس مبدأ، پورت مبدأ و طول payload تنظیم می شود. طول payload از مقدار ulen->udp استخراج شده و اندازه‌ی header UDP از آن کم می شود. داده‌ی payload با memmove داخل آرایه‌ی داخلى صفت کپی می شود و ایندکس نوشتمن زیاد می شود. در نهایت با wakeup()، فرایندهایی که در recv() منتظر داده بوده‌اند بیدار می شوند. پس از پایان این مراحل، buffer اولیه‌ی دریافت شده با kfree آزاد می شود تا از نشت حافظه جلوگیری شود.

فصل ۸

گام هفتم: مدیریت آزادسازی های buffer شبکه

در گام آخر نیز مدیریت کامل آزادسازی های buffer شبکه را بررسی و تکمیل کردیم تا از نشت حافظه و پر شدن تدریجی حافظه‌ی kernel جلوگیری شود.

در مسیر دریافت، اولیه‌ای که از تابع ip_rx() به net_rx_1000e() و سپس به ip_rx() منتقل می‌شود، در نهایت داخل ip_rx() آزاد می‌گردد. در صورتی که بسته به هر دلیل معتبر نباشد (مثل ناکافی بودن طول، غیر UDP بودن پروتکل، ثبت‌نشده بودن port مقصد یا پر بودن صفت مربوطه) همان‌جا با kfree(buf) آزاد می‌شود و دیگر وارد صفت نمی‌شود. این کار باعث می‌شود حتی در حالت drop شدن بسته‌ها نیز حافظه آزاد شود.

در حالتی که بسته معتبر باشد و در صفت مربوط به port ذخیره شود، داده‌ی payload داخل ساختار داخلی صفت کپی می‌شود و سپس buffer اصلی که از درایور دریافت شده بود آزاد می‌شود. در نتیجه صفت داخلی هر port فقط دیتای مورد نیاز را نگه می‌دارد و وابسته به buffer اولیه باقی نمی‌ماند و خب باعث می‌شود طول عمر buffer های تخصیص‌یافته توسط درایور کوتاه باشد و مدیریت حافظه ساده‌تر شود.

در سمت کاربر، هنگام اجرای recv()، پس از آنکه داده از صفت مربوط به port برداشته و به فضای کاربر کپی شد، صفت به روزرسانی می‌شود و آن خانه از صفت دوباره قابل استفاده می‌گردد. چون داده قبل‌کپی شده است، نیازی به نگهداشتن buffer اضافی وجود ندارد و هیچ pointer ای به حافظه آزادشده باقی نمی‌ماند.

در مسیر ارسال نیز آزادسازی buffer تنها پس از اطمینان از اتمام ارسال انجام می‌شود. همان‌طور که در تابع transmit_1000e() پیاده‌سازی شده، بعد از فعال شدن بیت وضعیت DD، buffer مربوط به آن descriptor با kfree آزاد می‌شود. این ترتیب باعث می‌شود حافظه پیش از اتمام استفاده‌ی سخت‌افزار آزاد نشود و در عین حال پس از پایان ارسال نیز بدون تأخیر آزاد گردد.

فصل ۹

سناریوی تست شبکه

برای تست نهایی مسیر شبکه، یک فایل TestScenario اضافه کردیم تا عملکرد UDP و محدودیت صف را بررسی کنیم. ابتدا در محیط xv6 روی یک port مشخص عملیات bind انجام می‌شود. سپس از سمت host و خارج از فولدر xv6، با استفاده از netcat چند پیام UDP به همان port ارسال می‌کنیم.

در خروجی داخل xv6 مشاهده می‌شود که بسته‌ها توسط ip_rx دریافت شده‌اند و داده‌ها از طریق recv() خوانده می‌شوند. در مرحله‌ی بعد، برای تست محدودیت صف، ۳۰ پیام UDP پشت سر هم ارسال می‌کنیم. طبق پیاده‌سازی، حداقل ۱۶ بسته در صف هر port نگه داشته می‌شود؛ بنابراین در خروجی دقیقاً ۱۶ پیام دریافت می‌شود و پیام‌های اضافی drop می‌شوند. پیام انتهايی نيز نشان می‌دهد که صف درست مدیریت شده و برنامه بدون crash اجرا را تمام کرده است.

(در ادامه تصاویر ترمینال host و ترمینال داخل xv6 آورده شده‌اند که ارسال بسته‌ها و دریافت ۱۶ پیام و drop شدن بقیه را نشان می‌دهند.)

A screenshot of a terminal window titled "Feb 14 2003". The terminal shows several commands being run:

```
user@user-VivoBook-ASUSLaptop-X513EQN-K513EQ:~/project$ echo fatimah | nc -u 127.0.0.1 26999
^C
user@user-VivoBook-ASUSLaptop-X513EQN-K513EQ:~/project$ echo hello | nc -u 127.0.0.1 26999
^C
user@user-VivoBook-ASUSLaptop-X513EQN-K513EQ:~/project$ for i in {1..30}; do echo pkt${i} | nc -u -w0 127.0.0.1 26999; done
user@user-VivoBook-ASUSLaptop-X513EQN-K513EQ:~/project$
```

شکل ۱-۹: ارسال بسته‌های host از سمت UDP برای تست دریافت و محدودیت صفت

A screenshot of a terminal window titled "Feb 14 2003". The terminal shows the output of a UDP stress test:

```
blind(10)
TX test (xv6 -> host)
RX test (host -> xv6). Now send UDP to host port 26999...
[...]
12:00:00.000000 167772674 got "fatimah"
" from src=167772674 sport=60691
queue/drop test: send 30 UDP packets quickly from host now...
receiving 16 packets (max queue size) ...
stress pkt 0: hello
stress pkt 1: pkt1
stress pkt 2: pkt2
stress pkt 3: pkt3
stress pkt 4: pkt4
stress pkt 5: pkt5
stress pkt 6: pkt6
stress pkt 7: pkt7
stress pkt 8: pkt8
stress pkt 9: pkt9
stress pkt 10: pkt10
stress pkt 11: pkt11
stress pkt 12: pkt12
stress pkt 13: pkt13
stress pkt 14: pkt14
stress pkt 15: pkt15
received 16 packets (should be 16). extra should have been dropped.
unbind(10)
fullnet: DONE (no crash, queue limit enforced, unbind drops packets)
```

شکل ۲-۹: خروجی پس از دریافت بسته‌های UDP

فصل ۱۰

بخش امتیازی: پشتیبانی از ICMP و گزارش خطای

در این بخش، برای بهبود رفتار پسته‌ی شبکه در مواجهه با خطاها ارتباطی، پشتیبانی از پیام‌های ICMP به سیستم اضافه شد. در پیاده‌سازی اولیه، اگر بسته‌ای به مقصد نمی‌رسید، فرایند کاربر از این موضوع مطلع نمی‌شد و ممکن بود در () recv یا هنگام ارسال، به صورت نامحدود منتظر بماند. با اضافه کردن پردازش پیام‌های ICMP، سیستم عامل می‌تواند خطاها مربوط به عدم دسترسی به مقصد را تشخیص داده و به فضای کاربر اعلام کند.

تغییرات اصلی در فایل kernel/net.c و داخل تابع net_rx انجام شد. در این تابع، علاوه بر بررسی بسته‌های ARP و IP/UDP، شرطی برای شناسایی بسته‌هایی با پروتکل IPPROTO_ICMP اضافه شد. در صورت تشخیص ICMP پردازش به تابع جداگانه‌ای هدایت می‌شود که سرآیند ICMP را تحلیل می‌کند.

در پردازش ICMP نوع پیام بررسی می‌شود و تمرکز اصلی روی پیام‌های نوع ۳ (Destination Unreachable) قرار دارد. این پیام شامل سرآیند IP و هشت بایت ابتدایی بسته‌ی اصلی است. از این اطلاعات، شماره پورت‌های مبدأ و مقصد UDP استخراج می‌شود تا مشخص شود کدام سوکت مسئول ارسال بسته‌ی اولیه بوده است. با استفاده از این اطلاعات، صفت مربوط به همان پورت در ساختار udp_queue شناسایی می‌شود.

برای اعلام خطا به فضای کاربر، ساختار udp_queue با یک فیلد وضعیت خطا تکمیل شد. هنگام دریافت پیام Destination Unreachable، این فیلد برای پورت مربوطه فعال می‌شود. در نتیجه، زمانی که فرایند کاربر تابع () recv یا () send را فراخوانی می‌کند، به جای انتظار بی‌پایان یا ارسال ناموفق، کد خطای مناسب دریافت می‌کند و از قطع ارتباط مطلع می‌شود.

با این تغییر، در صورت ارسال داده به یک میزبان غیرقابل دسترس، خطا در سطح kernel شناسایی شده و به صورت کنترل شده به برنامه‌ی کاربر منتقل می‌شود.