

# Résoudre les problèmes de fragmentation IP, MTU, MSS et PMTUD avec GRE et IPSEC

---

## Contenu

Introduction

### **Fragmentation et réassemblage IP**

Problèmes de fragmentation IP

Prévention de la fragmentation IP Tâches et fonctionnement de TCP MSS

### **Qu'est-ce que PMTUD ?**

Problèmes survenant avec la découverte PMTUD

Topologies réseau courantes nécessitant la découverte PMTUD

### **Qu'est-ce qu'un tunnel ?**

Considérations concernant les interfaces de tunnel

Le routeur possède un participant PMTUD au niveau du point d'extrémité d'un tunnel

"Mode Tunnel IPsec « pur »

### **GRE et IPsec ensemble**

**Autres recommandations**

**Informations connexes**

---

## Introduction

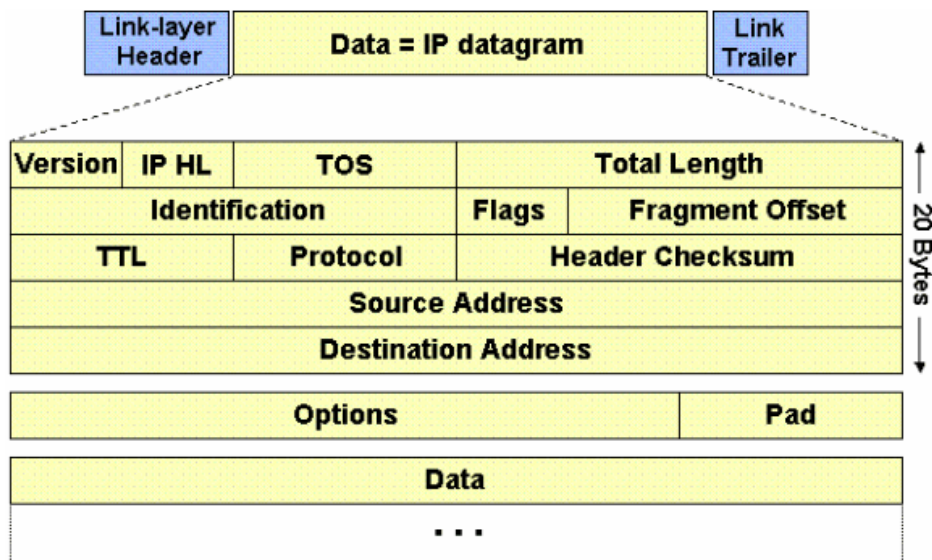
Le but de ce document est de présenter le fonctionnement de la fragmentation IP et de la découverte de l'unité de transmission maximale de la voie d'accès (PMTUD), ainsi que certains scénarios qui impliquent le comportement PMTUD associé à différentes combinaisons de tunnels IP. L'utilisation répandue des tunnels IP sur Internet a placé en première ligne les problèmes impliquant la fragmentation IP et la découverte PMTUD.

## Fragmentation et réassemblage IP

Le protocole IP a été conçu en vue d'une utilisation sur une grande variété de liaisons de transmission. Bien que la longueur maximale d'un datagramme IP soit 64K, la plupart des liaisons de transmission imposent une limite maximum de longueur de paquet inférieure, appelée MTU. La valeur de MTU dépend du type de liaison de transmission. La conception du protocole IP prend en charge les différences au niveau MTU en permettant à des routeurs de fragmenter les datagrammes IP, selon les besoins. La station de destination est chargée du réassemblage des fragments dans le datagramme IP initial (taille maximale).

La fragmentation IP implique de diviser un datagramme en un certain nombre de parties qui peuvent être réassemblées plus tard. La source IP, la destination, l'identification, la longueur totale et les zones de décalage de fragment, ainsi que les indicateurs « autres fragments » et « ne pas fragmenter » figurant dans l'en-tête IP, sont utilisés pour la fragmentation et le réassemblage IP. Pour plus d'informations sur les mécanismes de la fragmentation et du réassemblage IP, reportez-vous à RFC 791 . [↗](#)

L'image ci-dessous illustre la disposition d'un en-tête IP.



L'identification est 16 bits et est une valeur assignée par l'expéditeur d'un datagramme IP pour faciliter le réassemblage des fragments d'un datagramme.

Le décalage de fragment est 13 bits et indique l'emplacement d'un fragment dans le datagramme IP d'origine. Cette valeur est un multiple de huit octets.

Dans le champ d'indicateurs de l'en-tête IP, il y a trois bits pour les indicateurs de contrôle. Il est important de noter que le bit « ne pas fragmenter » (DF) joue un rôle central dans PMTUD, parce qu'il détermine si un paquet peut être fragmenté.

0 mordu est réservé, et est toujours placé à 0. a mordu 1 est le bit DF (0 = « peut fragmenter, » 1 = « ne font pas fragment »). Le bit 2 représente le bit MF (0 = « dernier fragment », 1 = « autres fragments »).

Valeur	Bit 0 réservé	Bit 1 DF	Bit 2 MF
0	0	Mai	Dernier
1	0	Ne faites pas	Plus

Le graphique ci-dessous montre un exemple de fragmentation. Si vous ajoutez toutes les longueurs de fragments IP, la valeur dépasse de 60 la longueur du datagramme IP d'origine. La raison pour laquelle la longueur totale est augmentée de 60 est la suivante : trois en-têtes IP supplémentaires ont été créées (une pour chaque fragment suivant le premier fragment).

Le premier fragment a un décalage de 0, la longueur de ce fragment est 1500 ; cela inclut 20 octets pour l'en-tête IP d'origine légèrement modifiée.

Le deuxième fragment a un décalage de 185 ( $185 \times 8 = 1480$ ), ce qui signifie que la partie "données" de ce fragment démarre 1480 octets dans le datagramme IP d'origine. La longueur de ce fragment est 1500 ; cela inclut l'en-tête IP supplémentaire créée pour ce fragment.

Le troisième fragment a un décalage de 370 ( $370 \times 8 = 2960$ ), ce qui signifie que la partie "données" de ce fragment démarre 2960 octets dans le datagramme IP d'origine. La longueur de ce fragment est 1500 ; cela inclut l'en-tête IP supplémentaire créée pour ce fragment.

Le quatrième fragment a un décalage de 555 ( $555 \times 8 = 4440$ ), ce qui signifie que la partie "données" de ce fragment démarre 4440 octets dans le datagramme IP d'origine. La longueur de ce fragment est de 700 octets ; cela inclut l'en-tête IP supplémentaire créée pour ce fragment.

La taille du datagramme IP d'origine ne peut être déterminée qu'une fois le dernier fragment reçu.

Le décalage de fragment du dernier fragment (555) entraîne un décalage de données de 4440 octets dans le datagramme IP d'origine. Si vous ajoutez ensuite les octets de données du dernier fragment ( $680 = 700 - 20$ ), on obtient 5120 octets, ce qui correspond à la partie "données" du datagramme IP d'origine. Ensuite, si vous ajoutez 20 octets pour une en-tête IP, vous obtenez la taille du datagramme IP d'origine ( $4440 + 680 + 20 = 5140$ ).

## Original IP Datagram

Sequence	Identifiant	Total Length	DF May / Don't	MF Last / More	Fragment Offset
0	345	5140	0	0	0

## IP Fragments (Ethernet)

Sequence	Identifiant	Total Length	DF May / Don't	MF Last / More	Fragment Offset
0-0	345	1500	0	1	0
0-1	345	1500	0	1	185
0-2	345	1500	0	1	370
0-3	345	700	0	0	555

## Problèmes de fragmentation IP

Il existe plusieurs problèmes qui rendent la fragmentation IP indésirable. On assiste à une faible augmentation de l'utilisation du processeur et de la mémoire pour la fragmentation d'un datagramme IP. Cela se vérifie pour l'expéditeur aussi bien que pour le routeur dans le chemin compris entre un expéditeur et un récepteur. La création de fragments implique simplement de créer des en-têtes et de copier le datagramme initial dans les fragments. Cela peut être fait assez efficacement, car toutes les informations requises pour la création des fragments est immédiatement disponible.

La fragmentation entraîne davantage de surcharge pour le récepteur lors du réassemblage des fragments, car le récepteur doit allouer de la mémoire pour les fragments en entrée et les fusionner de nouveau dans un datagramme une fois tous les fragments reçus. Le réassemblage sur un hôte n'est pas considéré comme un problème, parce que l'hôte possède les ressources temporelles et de mémoire suffisantes pour l'exécution de cette tâche.

Mais le réassemblage est totalement inefficace sur un routeur dont le rôle principal consiste à expédier des paquets le plus rapidement possible. Les routeurs ne sont pas conçus pour la rétention de paquets pendant une durée indéterminée. Par ailleurs, un routeur chargé d'un réassemblage choisit la plus grande mémoire tampon disponible (18K), car il ne peut connaître la taille du paquet IP initial qu'une fois le dernier fragment reçu.

Un autre problème de fragmentation réside dans la manipulation de fragments supprimés. Si un fragment de datagramme IP est supprimé, la totalité du datagramme IP d'origine doit être renvoyé, et il sera également fragmenté. Vous pouvez en voir un exemple avec le système NFS (Network File System). Par défaut, NFS a une taille de bloc de lecture/écriture de 8192 ; par conséquent, un datagramme NFS IP/UDP sera approximativement de 8500 octets (y compris les en-têtes NFS, UDP et IP). Une station émettrice connectée à Ethernet (MTU 1500) devra fragmenter le datagramme de 8500 octets en six parties : cinq fragments de 1500 octets et un fragment de 1100 octets. Si l'un des six fragments est supprimé en raison d'une liaison saturée, le datagramme initial complet devra être retransmis, ce qui signifie que six fragments supplémentaires devront être créés. Si cette liaison supprime l'un des six paquets, il y a peu de chances que les données NFS puissent être transférées sur cette liaison : en effet, un fragment IP minimum est dans ce cas supprimé de chaque datagramme IP d'origine (NFS 8500 octets).

Les pare-feu qui filtrent ou manipulent les paquets sur la base des informations de la couche 4 (L4) à la couche 7 (L7) peuvent rencontrer des difficultés de traitement des fragments IP. Si les fragments IP sont dans un ordre incorrect, un pare-feu peut bloquer les fragments non initiaux parce qu'ils ne diffusent pas les informations qui correspondent au filtre de paquet. Cela signifie que le datagramme IP d'origine ne peut pas être réassemblé par l'hôte de destination. Si le pare-feu est configuré de façon à autoriser les fragments non initiaux avec des informations insuffisantes en correspondance avec le filtre, une attaque de fragment non initial par le pare-feu risque de se produire. En outre, certains périphériques réseau (tels que les moteurs de contenus) acheminent les paquets sur la base des informations des couches L4 à L7, et si un paquet comporte des fragments multiples, le périphérique peut rencontrer des difficultés lors de la mise en oeuvre de ses stratégies.

## Prévention de la fragmentation IP Tâches et fonctionnement de TCP MSS

La taille maximum de segment TCP (MSS) définit la quantité maximale de données qu'un hôte souhaite accepter dans un datagramme TCP/IP simple. Ce datagramme TCP/IP peut être fragmenté au niveau de la couche IP. La valeur MSS est envoyée comme en-tête TCP uniquement dans les segments TCP SYN. Chaque côté d'une connexion TCP indique sa valeur MSS à l'autre côté. Contrairement à ce que l'on croit, la valeur MSS n'est pas négociée entre les hôtes. L'hôte expéditeur doit limiter la taille des données dans les segments TCP à une valeur inférieure ou égale au MSS indiqué par l'hôte de destination.

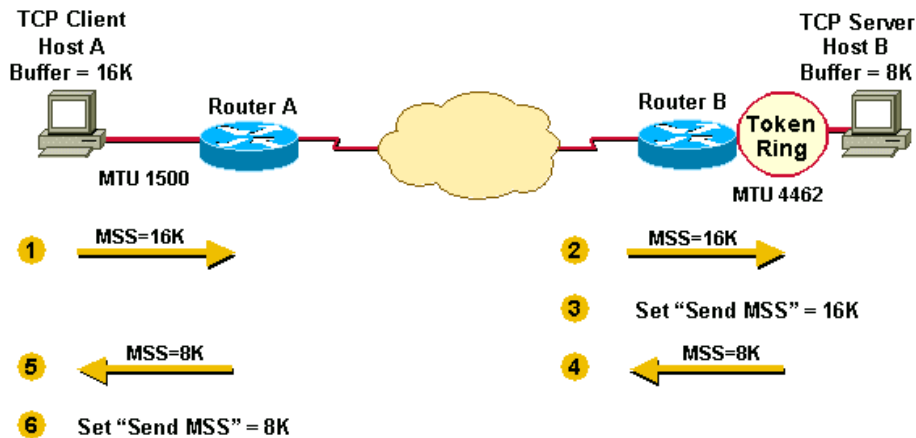
Initialement, la valeur MSS indiquait la taille d'une mémoire tampon (supérieure ou égale à 65496K) allouée à une station de destination, lui permettant de stocker les données TCP contenues dans un datagramme IP simple. MSS était le segment maximum de données (bloc) que le récepteur TCP souhaitait accepter. Ce segment TCP pouvait atteindre 64K (la taille maximum de datagramme IP) et pouvait être fragmenté au niveau de la couche IP afin d'être transmis via le réseau à l'hôte de destination. L'hôte de destination pouvait alors réassembler le datagramme IP avant de transmettre le segment TCP complet à la couche TCP.

Voici deux scénarios montrant comment les valeurs MSS sont définies et utilisées pour limiter les tailles de segments TCP, et donc les tailles de datagrammes IP.

Le scénario 1 illustre la façon dont MSS a été mis en application la première fois. L'hôte A a une mémoire tampon de 16K et l'hôte B une

mémoire tampon de 8K. Ils envoient et reçoivent leurs valeurs MSS et ajustent leur envoi MSS en vue de l'envoi réciproque de données. Notez que l'hôte A et l'hôte B devront fragmenter les datagrammes IP dont la taille dépasse celle de l'interface MTU mais reste inférieure à celle du MSS, car la pile TCP peut transmettre 16K ou 8K d'octets de données à destination d'IP. Dans le cas de l'hôte B, les paquets ont été fragmentés deux fois (une fois pour l'accès au réseau local Token Ring et une deuxième fois pour l'accès au réseau local Ethernet).

### Scénario 1



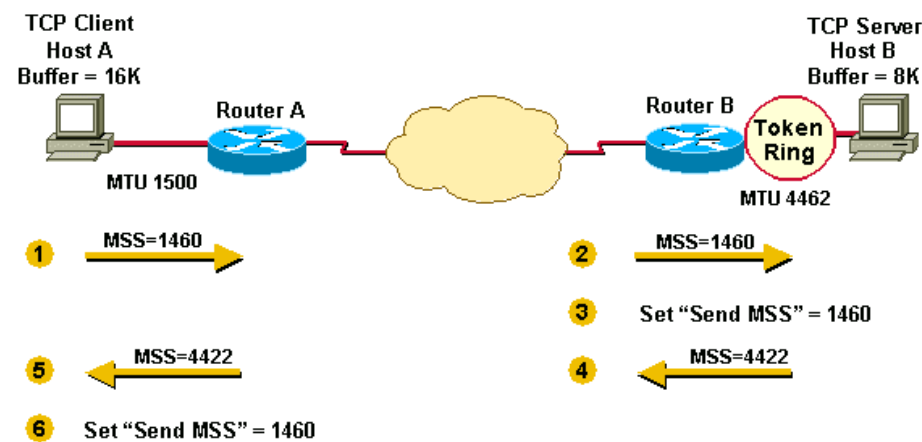
1. L'hôte A envoie sa valeur MSS de 16K à l'hôte B.
2. L'hôte B reçoit la valeur 16K MSS envoyée par l'hôte A.
3. L'hôte B définit sa valeur MSS à 16K.
4. L'hôte B envoie sa valeur MSS de 8K à l'hôte A.
5. L'hôte A reçoit la valeur 8K MSS envoyée par l'hôte A.
6. L'hôte A définit sa valeur MSS à 8K.

Afin d'éviter la fragmentation IP au niveau des points d'extrémité de la connexion TCP, la sélection de la valeur MSS a été modifiée ; elle est maintenant égale à la taille minimale de mémoire tampon et à la valeur MTU de l'interface sortante (- 40). Les valeurs MSS sont inférieures de 40 octets aux valeurs MTU, car MSS correspond à la taille de données TCP, ce qui n'inclut ni les 20 octets de l'en-tête IP, ni les 20 octets de l'en-tête TCP. MSS est basé sur les tailles d'en-tête par défaut ; la pile expéditeur doit soustraire les valeurs correspondant à l'en-tête IP et à l'en-tête TCP, selon les options TCP ou IP utilisées.

Le fonctionnement actuel de MSS est le suivant : chaque hôte compare tout d'abord son interface MTU sortante avec sa propre mémoire tampon, et choisit la valeur la plus faible en tant que valeur MSS à envoyer. Ensuite, les hôtes comparent la taille MSS reçue avec leur propre interface MTU, et choisissent à nouveau la plus faible des deux valeurs.

Le scénario 2 illustre cette étape supplémentaire exécutée par l'expéditeur pour éviter la fragmentation sur les fils locaux et distants. Vous remarquerez que la valeur MTU de l'interface sortante est prise en compte par chaque hôte (avant que les hôtes ne s'envoient mutuellement leurs valeurs MSS) et que cela permet d'éviter la fragmentation.

### Scénario 2



1. L'hôte A compare sa mémoire tampon MSS (16K) et son MTU (1500 - 40 = 1460) et utilise la valeur la plus basse comme MSS (1460) pour l'envoyer à l'hôte B.
2. L'hôte B reçoit la valeur MSS envoyée par l'hôte A (1460) et la compare à la valeur MTU de son interface sortante - 40 (4422).
3. L'hôte B définit la valeur la plus basse (1460) comme valeur MSS utilisée pour l'envoi des datagrammes IP à l'hôte A.
4. L'hôte B compare sa mémoire tampon MSS (8K) et son MTU (4462-40 = 4422), puis utilise 4422 comme MSS à envoyer à l'hôte A.
5. L'hôte A reçoit la valeur MSS envoyée par l'hôte B (4422) et la compare à la valeur de son interface de sortie MTU - 40 (1460).
6. L'hôte A définit la valeur la plus basse (1460) comme valeur MSS utilisée pour l'envoi des datagrammes IP à l'hôte B.

1460 est la valeur choisie par les deux hôtes comme valeur MSS à s'envoyer réciproquement. Souvent la valeur de l'envoi MSS sera identique

sur chaque extrémité d'une connexion TCP.

Dans le scénario 2, la fragmentation ne se produit pas au niveau des points d'extrémité d'une connexion TCP, car les deux valeurs MTU de l'interface sortante sont prises en compte par les hôtes. Les paquets peuvent toutefois être fragmentés au sein du réseau entre le routeur A et le routeur B, s'ils trouvent une liaison portant une valeur MTU inférieure à celle de l'interface sortante des hôtes.

## Qu'est-ce que PMTUD ?

Comme indiqué ci-dessus, TCP MSS effectue la fragmentation au niveau des deux points d'extrémité d'une connexion TCP, mais ne traite pas les cas de liaison MTU faible entre ces deux points. PMTUD a été développé pour éviter la fragmentation entre les points d'extrémité. Il est utilisé pour déterminer dynamiquement la valeur MTU la plus faible entre la source et la destination d'un paquet.

**Remarque:** PMTUD n'est pris en charge que par TCP. Le protocole UDP ainsi que d'autres protocoles ne le prennent pas en charge. Si PMTUD est activé sur un hôte (il l'est dans la grande majorité des cas), tous les paquets TCP/IP de l'hôte devront avoir le bit DF défini.

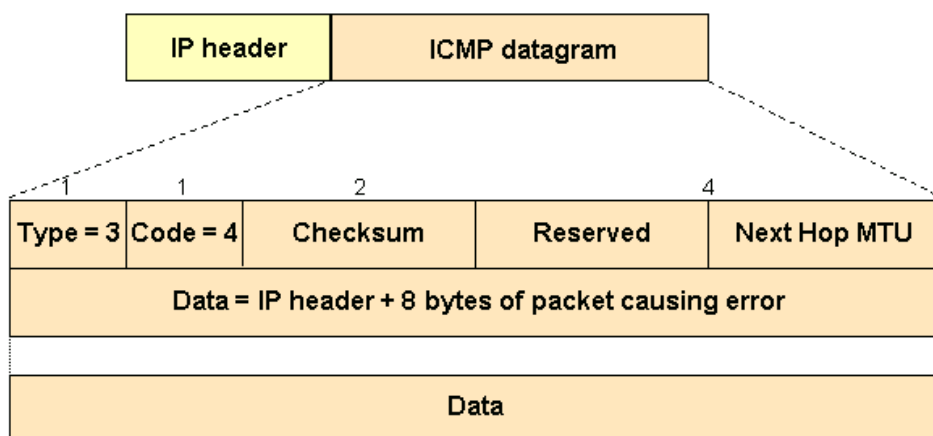
Lorsqu'un hôte envoie un paquet de données MSS complet avec bit DF défini, PMTUD réduit la valeur de l'envoi MSS pour la connexion s'il reçoit l'information indiquant que le paquet nécessite une fragmentation. Un hôte « se rappelle » généralement de la valeur MTU d'une destination en créant une entrée « hôte » (/32) dans sa table de routage, avec cette valeur MTU.

Si un routeur tente d'expédier un datagramme IP avec bit DF défini sur une liaison possédant une valeur MTU inférieure à la taille du paquet, ce routeur supprime le paquet et renvoie un message de destination inaccessible via le protocole ICMP (Internet Control Message Protocol) à la source de ce datagramme IP, avec le code indiquant « fragmentation requise et DF défini » (type 3, code 4). Lorsque la station source reçoit le message ICMP, elle diminue la valeur MSS ; quand TCP retransmet le segment, il utilisera la taille de segment la plus faible.

Voici un exemple de message ICMP « fragmentation requise et DF défini » susceptible de s'afficher sur un routeur après l'activation de la commande **debug ip icmp** :

```
ICMP: dst (10.10.10.10) frag. needed and DF set
unreachable sent to 10.1.1.1
```

Le diagramme ci-dessous montre le format de l'en-tête ICMP d'un message « fragmentation requise et DF défini » « destination inaccessible ».



Conformément à RFC 1191, un routeur qui renvoie un message ICMP indiquant « fragmentation requise et DF défini » doit inclure la valeur MTU de ce réseau dans les 16 bits du champ d'en-tête supplémentaire ICMP étiqueté « inutilisée » dans la spécification ICMP RFC 792 de . [☞](#)

Les premières mises en oeuvre de RFC 1191 n'ont pas fourni les informations MTU de saut suivant. Même lorsque ces informations ont été fournies, certains hôtes les ignorent. Pour ces situations, RFC 1191 contient également un tableau qui présente les valeurs suggérées par lesquelles la valeur MTU doit être diminuée pendant le PMTUD. Il est utilisé par les hôtes pour parvenir plus rapidement à une valeur raisonnable pour l'envoi MSS.

Plateau	MTU	Comments	Reference
-----	---	-----	-----
	65535	Official maximum MTU	RFC 791
	65535	Hyperchannel	RFC 1044
65535			
32000		Just in case	
	17914	16Mb IBM Token Ring	ref. [6]
17914			
	8166	IEEE 802.4	RFC 1042
8166			
	4464	IEEE 802.5 (4Mb max)	RFC 1042
	4352	FDDI (Revised)	RFC 1188
4352 (1%)			
	2048	Wideband Network	RFC 907
	2002	IEEE 802.5 (4Mb recommended)	RFC 1042
2002 (2%)			
	1536	Exp. Ethernet Nets	RFC 895
	1500	Ethernet Networks	RFC 894
	1500	Point-to-Point (default)	RFC 1134
	1492	IEEE 802.3	RFC 1042
1492 (3%)			
	1006	SLIP	RFC 1055
	1006	ARPANET	BBN 1822
1006			
	576	X.25 Networks	RFC 877
	544	DEC IP Portal	ref. [10]
	512	NETBIOS	RFC 1088
	508	IEEE 802/Source-Rt Bridge	RFC 1042
	508	ARCNET	RFC 1051
508 (13%)			
	296	Point-to-Point (low delay)	RFC 1144
296			
68		Official minimum MTU	RFC 791

PMTUD est exécuté continuellement sur tous les paquets, car le chemin entre l'expéditeur et le destinataire peut changer dynamiquement. Chaque fois qu'un expéditeur reçoit un message ICMP « fragmentation impossible », il met à jour les informations de routage (dans lesquelles il stocke le PMTUD).

Deux situations peuvent se produire pendant le PMTUD :

- Le paquet peut être acheminé jusqu'au destinataire sans être fragmenté.

**Remarque:** Pour qu'un routeur puisse protéger le processeur contre les attaques DoS, il limite le nombre de messages ICMP qu'il envoie à deux par seconde. Par conséquent, dans ce contexte, si vous avez un scénario réseau dans lequel vous prévoyez que le routeur doit répondre à plus de deux ICMP (code = 3, type = 4) par seconde (sur des hôtes différents le cas échéant), vous pouvez désactiver l'étranglement des messages ICMP à l'aide de la commande **aucun ip icmp rate-limit unreachable [df] interface** .

- L'expéditeur peut recevoir des messages ICMP « fragmentation impossible » en provenance de sauts situés sur le chemin du destinataire.

PMTUD est exécuté indépendamment pour les deux directions d'un flux TCP. Dans certains cas, PMTUD déclenche l'une des stations dans une direction de flux pour diminuer le MSS, tandis que l'autre station conserve le MSS d'origine, car elle n'a jamais envoyé de datagramme IP suffisamment volumineux pour justifier un déclenchement de PMTUD.

Vous trouverez un exemple de cette situation dans le scénario 3 (connexion HTTP). Le client TCP envoie des paquets de faible taille et le serveur envoie des paquets volumineux. Dans ce cas, seuls les paquets volumineux (supérieurs à 576 octets) déclenchent PMTUD. Les paquets peu volumineux (inférieurs à 576 octets) ne déclenchent pas PMTUD, car ils n'exigent pas de fragmentation sur la liaison MTU 576.

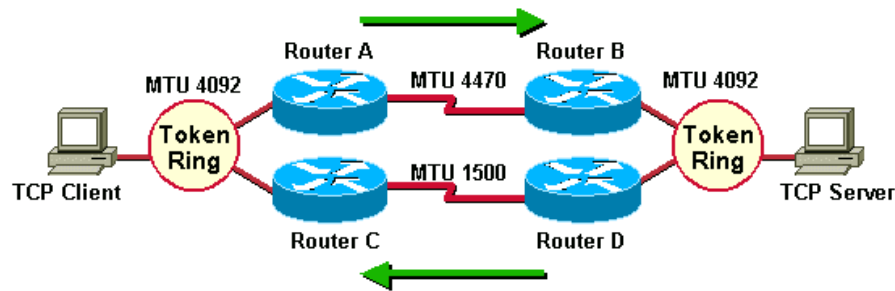
### Scénario 3



Le scénario 4 montre un exemple de routage asymétrique, dans lequel l'un des chemins affiche une valeur MTU minimale inférieure à celle de l'autre. Un routage asymétrique se produit lorsque différents chemins sont empruntés pour l'envoi et la réception de données entre deux points d'extrémité. Dans ce scénario, PMTUD déclenche la diminution de la valeur MSS uniquement dans une direction de flux TCP. Le trafic entre le client TCP et le serveur passe par le routeur A et le routeur B, tandis que le trafic retour (du serveur au client) passe par le routeur D et le routeur C. Quand le serveur TCP envoie des paquets au client, PMTUD déclenche le serveur afin de diminuer la valeur MSS, car le routeur D doit fragmenter les paquets de 4092 octets avant de pouvoir les envoyer au routeur C.

Le client, quant à lui, ne reçoit pas de message ICMP « destination inaccessible » avec le code indiquant « fragmentation requise et DF défini », car le routeur A ne doit pas effectuer la fragmentation des paquets lors de l'envoi au serveur via le routeur B.

## Scénario 4



**Remarque:** La commande `ip tcp path-mtu-discovery` est utilisée pour activer la découverte des voies d'accès MTU TCP pour les connexions TCP initialisées par des routeurs (BGP et Telnet par exemple).

## Problèmes survenant avec la découverte PMTUD

Il existe trois problèmes potentiels au niveau de PMTUD, dont deux rares et un fréquent.

- Un routeur peut supprimer un paquet et ne pas envoyer de message ICMP. (Rare)
- Un routeur peut générer et envoyer un message ICMP, mais ce dernier peut être bloqué par un routeur ou par un pare-feu situé entre ce routeur et l'expéditeur. (Terrain communal)
- Un routeur peut générer et envoyer un message ICMP, mais l'expéditeur peut ignorer ce message. (Rare)

Le premier et le dernier point de la liste ci-dessous correspondent à des situations rares et proviennent généralement d'une erreur, tandis que le point intermédiaire décrit un problème fréquent. Les personnes chargées de la mise en oeuvre de filtres de paquets ICMP ont tendance à bloquer tous les types de messages ICMP, au lieu de ne bloquer que certains d'entre eux. Un filtre de paquets peut bloquer tous les types de messages ICMP à l'exception de ceux qui indiquent « inaccessible » ou « dépassement de délai ». Le succès ou l'échec de PMTUD est lié aux messages d'inaccessibilité ICMP parvenant à l'expéditeur d'un paquet TCP/IP. Les messages ICMP time-exceeded sont importants pour d'autres problèmes d'IP. Un exemple de filtre mis en oeuvre sur un routeur figure ci-dessous.

```
access-list 101 permit icmp any any unreachable
access-list 101 permit icmp any any time-exceeded
access-list 101 deny icmp any any
access-list 101 permit ip any any
```

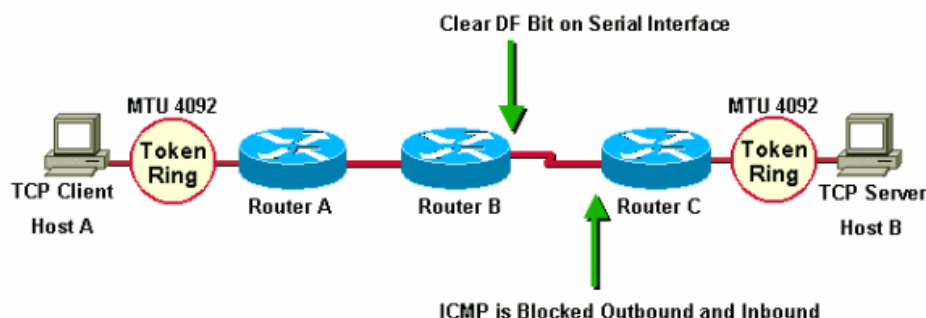
Il y a d'autres techniques qui peuvent être utilisées pour alléger le problème de blocage total de l'ICMP.

- Vous pouvez effacer le bit DF sur le routeur et autoriser malgré tout la fragmentation (ceci peut toutefois ne pas être une bonne idée). Pour plus d'informations, voir Problèmes de fragmentation IP.
- Manipulez la valeur de l'option MSS TCP à l'aide de la commande d'interface `ip tcp adjust-mss le <500-1460>`.

Dans le scénario 5 ci-dessous, le routeur A et le routeur B sont dans le même domaine administratif. Le routeur C est inaccessible et bloque ICMP ; par conséquent, PMTUD est interrompu. Une solution consisterait à effacer le bit DF dans les deux sens sur le routeur B, pour permettre la fragmentation. Ceci peut être fait en utilisant le routage spécifique. La syntaxe à utiliser pour effacer le bit DF est disponible dans Cisco IOS® 12.1(6) et versions ultérieures.

```
interface serial0
...
ip policy route-map clear-df-bit
route-map clear-df-bit permit 10
match ip address 111
set ip df 0

access-list 111 permit tcp any any
```



Une autre possibilité consisterait à changer la valeur de l'option MSS TCP sur les paquets SYN qui traversent le routeur (disponible dans Cisco

IOS 12.2(4)T et versions ultérieures). Ceci réduit la valeur de l'option MSS dans le paquet SYN TCP, de sorte qu'il soit plus petit que la valeur (1460) figurant dans la commande **ip tcp adjust-mss**. Le résultat est que l'expéditeur TCP enverra des segments qui ne dépasseront pas cette valeur. La taille de paquets IP sera 40 octets plus grande (1500) que la valeur MSS (1460 octets) pour tenir compte de l'en-tête TCP (20 octets) et de l'en-tête IP (20 octets).

Vous pouvez ajuster le MSS des paquets SYN TCP à l'aide de la commande **ip tcp adjust-mss**. La syntaxe suivante ramènera la valeur MSS des segments de TCP à 1460. Cette commande permet de traiter le trafic entrant et sortant sur l'interface serial0.

```
int s0
ip tcp adjust-mss 1460
```

Les problèmes de fragmentation IP sont devenus plus répandus depuis le déploiement massif des tunnels IP. La raison pour laquelle les tunnels entraînent davantage de fragmentation est liée au fait que l'encapsulation ajoute du « temps système » à la taille d'un paquet. Par exemple, l'encapsulation GRE (Generic Router Encapsulation) ajoute 24 octets à un paquet ; après cette augmentation, le paquet peut avoir besoin d'être fragmenté car il dépasse la taille de l'interface de sortie MTU. Dans une section ultérieure de ce document, vous verrez des exemples de problèmes qui peuvent se poser avec les tunnels et la fragmentation IP.

## Topologies réseau courantes nécessitant la découverte PMTUD

PMTUD est nécessaire dans les situations réseau dans lesquelles des liaisons intermédiaires possèdent des MTU plus faibles que le MTU des liaisons finales. Les raisons de ces valeurs plus faibles de MTU peuvent notamment être les suivantes :

- eToken Ring (ou FDDI) - hôtes d'extrémité connectés avec une connexion Ethernet entre elles. Les MTU Token Ring (ou FDDI) aux extrémités sont supérieurs au MTU Ethernet du milieu.
- PPPoE (souvent utilisé avec ADSL) a besoin de 8 octets pour son en-tête. Ceci ramène le MTU efficace d'Ethernet à 1492 (1500 - 8).

Les protocoles de canalisation en tunnel comme GRE, IPsec, et L2TP ont besoin également d'espace pour leurs en-têtes et queues de bande respectifs. Cela réduit également le MTU efficace de l'interface sortante.

Dans les sections suivantes nous étudierons l'incidence de PMTUD lorsqu'un protocole de canalisation en tunnel est utilisé entre les deux hôtes d'extrémité. Sur les trois cas ci-dessus, ce cas est le plus complexe, et il couvre tous les problèmes que vous pourriez rencontrer dans les autres cas.

## Qu'est-ce qu'un tunnel ?

Un tunnel est une interface logique sur un routeur Cisco qui fournit une méthode d'encapsulation de paquets passagers au sein d'un protocole de transport. C'est une architecture conçue pour fournir les services destinés à mettre en application un plan d'encapsulation point à point. La transmission tunnel possède les trois composants primaires suivants :

- Protocole passager (AppleTalk, Banyan VINES, CLNS, DECNet, IP, ou IPX)
- Protocole transporteur - L'un des protocoles d'encapsulation suivants :
  - GRE - Le protocole transporteur multiprotocole de Cisco. Voir RFC 2784 et RFC 1701 pour plus d'informations. [↗](#) [↗](#)
  - IP dans des tunnels IP - Voir RFC 2003 pour plus d'informations. [↗](#)
- Protocole de transport - Protocole utilisé pour acheminer le protocole encapsulé

Les paquets ci-dessous illustrent les concepts de transmission tunnel IP dans lesquels GRE est le protocole d'encapsulation et IP est le protocole de transport. Le protocole passager est également IP. Dans ce cas, IP représente à la fois le protocole de transport et le protocole passager.

### Paquet normal

IP	TCP	Telnet
----	-----	--------

### Paquet tunnel

IP	GRE	IP	TCP	Telnet
----	-----	----	-----	--------

•

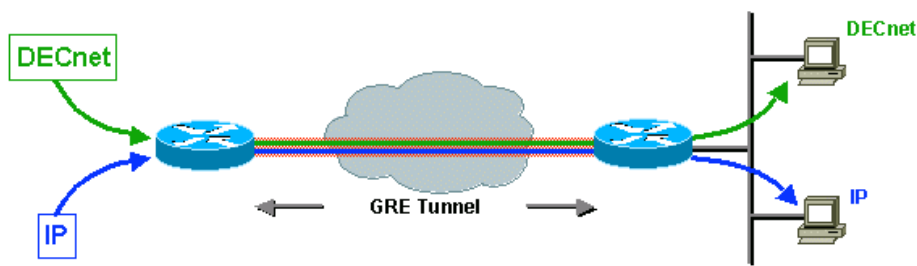
IP est le protocole de transport.

GRE est le protocole d'encapsulation.

IP est le protocole passager.

L'exemple suivant montre l'encapsulation d'IP et de DECNet comme protocoles passagers avec GRE comme transporteur. Cela illustre le fait que le protocole transporteur peut encapsuler des protocoles passagers multiples.





Un administrateur réseau peut envisager la transmission tunnel dans une situation où deux réseaux non contigus et non IP sont séparés par un circuit principal IP. Si les réseaux non contigus exécutent DECNet, l'administrateur peut ne pas vouloir les connecter ensemble en configurant DECNet dans le circuit principal. L'administrateur peut ne pas vouloir laisser le routage DECNet consommer la bande passante du réseau principal parce que cela pourrait affecter les performances du réseau IP.

Une alternative viable est la transmission tunnel de DECNet sur le circuit principal IP. La transmission tunnel permet d'encapsuler les paquets DECNet au sein du protocole IP et de les envoyer via le circuit principal au point d'extrémité du tunnel, où l'encapsulation est supprimée et où les paquets DECNet peuvent être routés vers leur destination au moyen du protocole DECNet.

Le trafic à l'intérieur d'un autre protocole fournit les avantages suivants :

- Les points d'extrémité utilisent des adresses privées (RFC 1918 ) et le circuit principal ne supporte pas le routage de ces adresses. [☞](#)
- Autorisez les VPN (Virtual Private Network) sur les WAN ou sur Internet.
- Joignez les réseaux multiprotocole non contigus sur un circuit principal à protocole unique.
- Cryptez le trafic sur le circuit principal ou sur Internet.

Pour le reste du document, nous utiliserons IP comme protocole passager et IP comme protocole de transport.

## Considérations concernant les interfaces de tunnel

Les informations suivantes sont des considérations liées à la transmission tunnel.

- La commutation rapide des tunnels GRE a été introduite dans le logiciel Cisco IOS version 11.1 et la commutation CEF a été introduite dans la version 12.0. La commutation CEF pour les tunnels GRE multipoints a été introduite dans la version 12.2(8)T. L'encapsulation et la suppression de l'encapsulation aux extrémités du tunnel étaient des opérations lentes dans les versions antérieures d'IOS, quand seule la commutation de processus était prise en charge.
- Lors de la transmission tunnel de paquets, des problèmes de sécurité et de topologie se posent. Les tunnels peuvent ignorer les listes de contrôle d'accès (ACL) et les pare-feu. Si vous effectuez une transmission tunnel via un pare-feu, vous évitez le pare-feu pour tout protocole passager impliqué dans la transmission tunnel. Par conséquent, il est recommandé d'inclure la fonction de pare-feu aux extrémités du tunnel, afin de mettre en oeuvre une politique au niveau des protocoles passagers.
- La transmission tunnel peut créer des problèmes avec les protocoles de transport, qui ont des temporisations limitées (DECnet, par exemple) en raison d'une latence accrue
- La transmission tunnel au sein d'environnements avec liaisons à vitesses variables (tels que les anneaux FDDI) et via des lignes téléphoniques 9600-bps lentes peut entraîner des problèmes de réorganisation des paquets. Certains protocoles passagers fonctionnent mal dans des réseaux de supports mixtes.
- Les tunnels point à point peuvent épuiser la bande passante sur une liaison physique. Si vous utilisez des protocoles de routage sur des tunnels point à point multiples, souvenez-vous que chaque interface de tunnel a une bande passante et que l'interface physique sur laquelle le tunnel est exécuté a elle-même une bande passante. Par exemple, vous pouvez définir la bande passante du tunnel à 100 Kb s'il y a 100 tunnels sur une liaison de 10 Mb. La bande passante par défaut pour un tunnel est 9Kb.
- Les protocoles de routage peuvent préférer a tunnel sur une liaison « réelle », parce que le tunnel pourrait sembler être une liaison à un saut avec le coût le plus faible, alors qu'en réalité il implique davantage de sauts et il est plus coûteux qu'un autre chemin. Cela peut être atténué avec la configuration appropriée du protocole de routage. Vous pouvez envisager d'exécuter un protocole de routage différent sur l'interface de tunnel que celui exécuté sur l'interface physique.
- Les problèmes de routage récursif peuvent être évités via la configuration appropriée de routes statiques vers la destination du tunnel. On parle de route récursif quand le meilleur chemin vers la destination du tunnel récursif est le tunnel lui-même. Cette situation entraîne un rebond de l'interface du tunnel. Vous pouvez voir l'erreur suivante en cas de problème de routage récursif.

```
%TUN-RECURDOWN Interface Tunnel 0
temporarily disabled due to recursive routing
```

## Le routeur possède un participant PMTUD au niveau du point d'extrémité d'un tunnel

Le routeur a deux rôles PMTUD différents à jouer lorsqu'il représente le point d'extrémité d'un tunnel.

- Dans le premier rôle, le routeur est l'expéditeur d'un paquet hôte. Pour le traitement PMTUD, le routeur doit contrôler le bit DF et la taille de paquet du paquet de données d'origine, puis exécuter l'action appropriée, le cas échéant.
- Le deuxième rôle intervient une fois que le routeur a encapsulé le paquet IP d'origine dans le paquet de tunnel. À ce stade, le routeur agit plutôt comme un hôte par rapport à PMTUD et par rapport au paquet IP du tunnel.

Examinons tout d'abord ce qui se passe lorsque le routeur agit dans le cadre du premier rôle (un routeur acheminant des paquets IP hôte) au niveau de PMTUD. Ce rôle intervient avant l'encapsulation du paquet IP hôte par le routeur dans le paquet du tunnel.

Si le routeur participe en tant que transporteur d'un paquet hôte, il doit exécuter les opérations suivantes :

- Vérifier si le bit DF a été défini.
- Vérifier la taille de paquet utilisable par le tunnel.
- Fragmenter (si le paquet est trop volumineux et si le bit DF n'est pas défini), encapsuler et envoyer les fragments ; ou
- Supprimer le paquet (si le paquet est trop volumineux et si le bit DF est défini) et envoyer un message ICMP à l'expéditeur.
- Effectuer l'encapsulation (si le paquet n'est pas trop volumineux) et l'envoi.

Vous pouvez choisir entre une encapsulation suivie d'une fragmentation (envoi de deux fragments d'encapsulation) et une fragmentation suivie d'une encapsulation (envoi de deux fragments encapsulés).

Vous trouverez ci-après quelques exemples du mécanisme d'encapsulation et de fragmentation de paquets IP, ainsi que deux scénarios qui montrent l'interaction de PMTUD et des paquets qui traversent les exemples de réseaux.

Le premier exemple ci-dessous montre ce qui arrive à un paquet quand le routeur (à la source du tunnel) joue le rôle de routeur transporteur. Souvenez-vous que pour le traitement PMTUD, le routeur doit contrôler le bit DF et la taille de paquet du paquet de données d'origine, puis exécuter l'action appropriée, le cas échéant. Ces exemples utilisent l'encapsulation GRE pour le tunnel. Comme vous pouvez le voir ci-dessous, GRE effectue la fragmentation avant l'encapsulation. Les exemples suivants montrent les scénarios dans lesquels la fragmentation est faite après l'encapsulation.

Dans l'exemple 1, le bit DF n'est pas défini ( $DF = 0$ ) et le tunnel GRE IP MTU est 1476 ( $1500 - 24$ ).

### Exemple 1

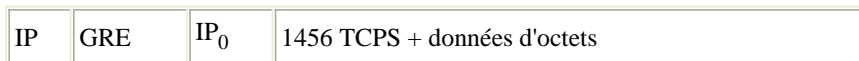
1. Le routeur transporteur (à la source du tunnel) reçoit un datagramme de 1500 octets avec bit DF envoyé ( $DF = 0$ ) en provenance de l'hôte expéditeur. Ce datagramme se compose d'un en-tête IP de 20 octets et d'une charge utile TCP de 1480 octets.



2. Puisque le paquet sera trop volumineux pour le MTU IP après ajout du temps système GRE (24 octets), le routeur transporteur divise le datagramme en deux fragments de 1476 (en-tête IP de 20 octets + charge utile IP de 1456 octets) plus 44 octets (20 octets d'en-tête IP + 24 octets de charge utile IP) ; par conséquent, une fois l'encapsulation GRE ajoutée, le paquet n'est pas plus volumineux que le MTU de l'interface physique sortante.



3. Le routeur transporteur ajoute l'encapsulation GRE (qui inclut un en-tête GRE de 4 octets et un en-tête IP de 20 octets) à chaque fragment du datagramme IP d'origine. Ces deux datagrammes IP ont maintenant une longueur de 1500 et 68 octets, et ces datagrammes sont vus en tant que datagrammes IP et non en tant que fragments.



4. Le tunnel de destination retire l'encapsulation GRE de chaque fragment du datagramme initial, en laissant deux fragments IP d'une longueur de 1476 et 24 octets. Ces fragments de datagrammes IP seront expédiés séparément par ce routeur au hôte de destination.



5. L'hôte de destination effectuera le réassemblage de ces deux fragments pour reconstituer le datagramme d'origine.



Le scénario 5 dépeint le rôle du routeur transporteur dans le contexte d'une topologie réseau.

Dans l'exemple suivant, le routeur joue le même rôle que le routeur transporteur, mais cette fois le bit DF est défini (DF = 1).

### Exemple 2

1. Le routeur transporteur (à la source du tunnel) reçoit un datagramme de 1500 octets avec bit DF= 1) en provenance de l'hôte expéditeur.



2. Puisque le bit DF est défini et que la taille du datagramme (1500 octets) est plus grande que le MTU IP du tunnel GRE (1476), le routeur supprimera le datagramme et enverra un message à la source du datagramme, indiquant que la fragmentation ICMP est requise mais que le bit DF est défini. Le message ICMP alertera l'expéditeur que le MTU est 1476.



3. L'hôte expéditeur reçoit le message ICMP, et quand il renvoie les données d'origine, il utilise un datagramme IP de 1476 octets.



4. Cette longueur de datagramme IP (1476 octets) est maintenant égale à la valeur du MTU IP du tunnel GRE. Par conséquent, le routeur ajoute l'encapsulation GRE au datagramme IP.



5. Le routeur de destination (à la destination du tunnel) supprime l'encapsulation GRE du datagramme IP et l'envoie à l'hôte destinataire.



Nous pouvons maintenant examiner ce qui se produit quand le routeur joue le deuxième rôle en tant qu'hôte expéditeur, par rapport au PMTUD et au paquet IP du tunnel. Rappelez-vous que ce rôle intervient une fois que le routeur a encapsulé le paquet IP d'origine dans le paquet de tunnel.

**Remarque:** Par défaut, un routeur n'utilise pas PMTUD sur les paquets du tunnel GRE qu'il génère. La commande **tunnel path-mtu-discovery** commande peut être utilisée pour activer PMTUD pour les paquets de tunnel GRE-IP.

Vous trouverez ci-après un exemple de ce qui se produit quand l'hôte envoie des datagrammes IP assez petits pour être inclus dans le MTU IP de l'interface de tunnel GRE. Dans ce cas, le bit DF peut être soit défini, soit effacé (1 ou 0). L'interface de tunnel GRE n'a pas la commande **tunnel path-mtu-discovery** configurée ; par conséquent, le routeur n'utilisera pas PMTUD sur le paquet GRE-IP.

### Exemple 3

1. Le routeur transporteur (à la source du tunnel) reçoit un datagramme de 1476 octets en provenance de l'hôte expéditeur.



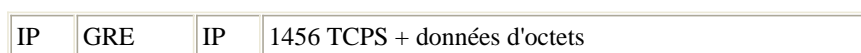
2. Ce routeur encapsule le datagramme IP de 1476 octets à l'intérieur de GRE afin d'obtenir un datagramme GRE IP de 1500 octets. Le bit DF de l'en-tête GRE IP est effacé (DF = 0). Ce routeur achemine ensuite ce paquet vers la destination du tunnel.



3. Supposons qu'il y ait un routeur entre la source et la destination du tunnel, avec un MTU de liaison de 1400. Ce routeur fragmentera le paquet du tunnel, car le bit DF est effacé (DF = 0). Rappelez-vous que dans cet exemple, l'IP extérieur est fragmenté, ainsi que le GRE, et que les en-têtes TCP ne sont affichées que dans le premier fragment.



4. Le routeur de destination du tunnel doit rassembler le paquet de tunnel GRE.



5. Une fois le paquet de tunnel GRE réassemblé, le routeur retire l'en-tête IP GRE et envoie le datagramme IP d'origine à sa destination.





L'exemple suivant montre ce qui se produit quand le routeur joue le rôle d'un hôte expéditeur par rapport au PMTUD et au paquet IP du tunnel. Cette fois, le bit DF est défini (DF = 1) dans l'en-tête IP initial et nous avons configuré la commande **tunnel path-mtu-discovery** de sorte que le bit DF soit copié de l'en-tête IP interne à l'en-tête externe (GRE + IP).

#### Exemple 4

1. Le routeur transporteur (à la source du tunnel) reçoit un datagramme de 1476 octets avec bit DF= 1) en provenance de l'hôte expéditeur.



2. Ce routeur encapsule le datagramme IP de 1476 octets à l'intérieur de GRE afin d'obtenir un datagramme GRE IP de 1500 octets. Cette en-tête IP GRE porte le bit DF défini (DF = 1) puisque le datagramme IP d'origine le portait également. Ce routeur achemine ensuite ce paquet vers la destination du tunnel.



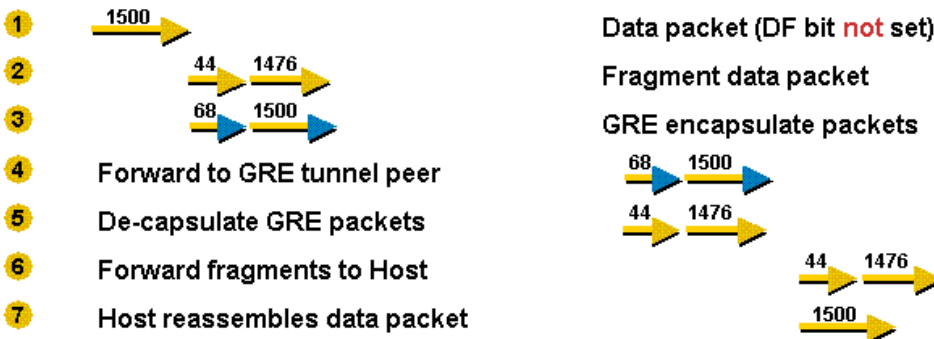
3. Supposons de nouveau qu'il y ait un routeur entre la source et la destination du tunnel, avec un MTU de liaison de 1400. Ce routeur ne fragmentera pas le paquet de tunnel, car le bit DF est défini (DF = 1). Ce routeur doit supprimer le paquet et envoyer un message d'erreur ICMP au routeur source du tunnel, puisque c'est l'adresse IP de la source sur le paquet.



4. Le routeur transporteur (à la source du tunnel) reçoit ce message d'erreur ICMP et diminue le MTU IP du tunnel GRE en faisant passer sa valeur à 1376 (1400 - 24). La prochaine fois que l'hôte expéditeur retransmettra les données au sein d'un paquet IP de 1476 octets, ce paquet sera trop volumineux et ce routeur enverra un message d'erreur ICMP à l'expéditeur, portant la valeur MTU de 1376. Quand l'hôte expéditeur retransmet les données, il envoie un paquet IP de 1376 octets, et ce paquet est acheminé via le tunnel GRE jusqu'à l'hôte destinataire.

#### Scénario 5

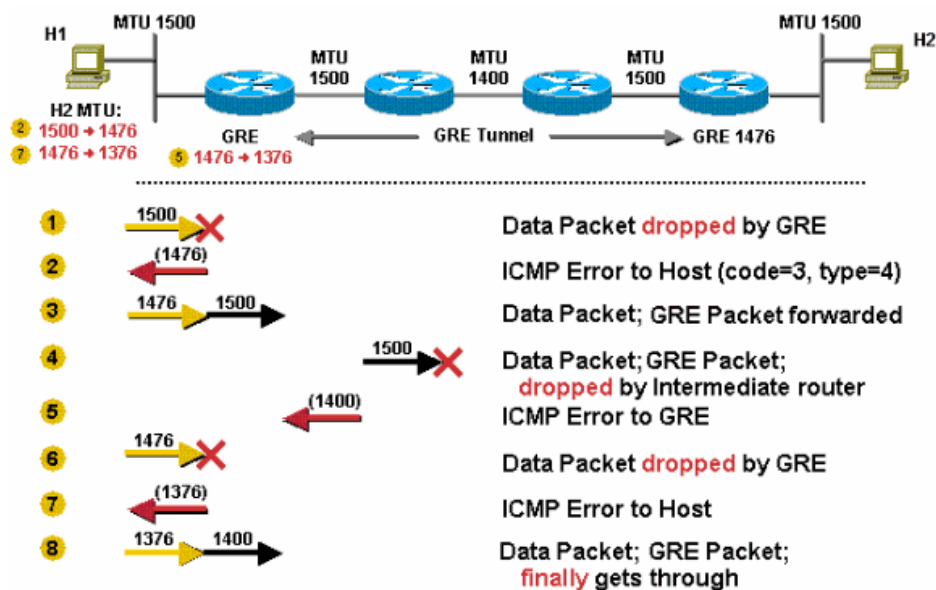
Ce scénario illustre la fragmentation GRE. Rappelez-vous que vous devez fragmenter avant l'encapsulation pour GRE, puis que vous devez appliquer PMTUD pour le paquet de données ; notez également que le bit DF n'est pas copié lorsque le paquet IP est encapsulé par GRE. Dans ce scénario, le bit DF n'est pas défini. Le MPU IP de l'interface de tunnel GRE est par défaut de 24 octets, soit inférieur au MTU IP de l'interface physique. Le MTU IP de l'interface GRE est donc de 1476.



1. L'expéditeur envoie un paquet de 1500 octets (en-tête IP de 20 octets + charge utile TCP de 1480 octets).
2. Puisque le MTU du tunnel GRE est 1476, le paquet de 1500 octets est divisé en deux fragments IP de 1476 et 44 octets, chacun en prévision des 24 octets supplémentaires de l'en-tête GRE.
3. Les 24 octets de l'en-tête GRE sont ajoutés à chaque fragment IP. Maintenant les fragments sont de 1500 (1476 + 24) et 68 (44 + 24) octets chacun.
4. Les paquets GRE + IP contenant les deux fragments IP sont expédiés au routeur partenaire de tunnel GRE.
5. Le routeur partenaire de tunnel GRE retire les en-têtes GRE des deux paquets.
6. Ce routeur achemine les deux paquets vers l'hôte de destination.
7. L'hôte de destination réassemble les fragments IP dans le datagramme IP d'origine.

## Scénario 6

Ce scénario est similaire au scénario 5, mais cette fois-ci, le bit DF est défini. Dans le scénario 6, le routeur est configuré de façon à appliquer PMTUD sur GRE + sur les paquets IP de tunnel à l'aide de la commande **tunnel path-mtu-discovery**, et le bit DF est copié de l'en-tête d'origine vers l'en-tête IP GRE. Si le routeur reçoit une erreur ICMP pour le GRE + paquet IP, il diminue la valeur MTU IP sur l'interface de tunnel GRE. De nouveau, rappelez-vous que le tunnel MTU IP du tunnel GRE est égal à 24 octets, soit moins que le MTU de l'interface physique par défaut ; par conséquent, dans ce cas le MTU IP GRE est égal à 1476. Notez également qu'il y a une liaison MTU 1400 dans le chemin du tunnel GRE.



1. Le routeur reçoit un paquet de 1500 octets (en-tête IP de 20 octets + charge utile TCP de 1480), et il supprime le paquet. Le routeur supprime le paquet, car il est plus volumineux que le MTU IP (1476) de l'interface de tunnel GRE.
2. Le routeur envoie une erreur ICMP à l'expéditeur lui indiquant que le MTU du saut suivant est de 1476. L'hôte enregistre ces informations, habituellement sous forme de route hôte de destination dans sa table de routage.
3. L'hôte expéditeur utilise une taille de paquet de 1476 octets lorsqu'il envoie de nouveau les données. Le routeur GRE ajoute 24 octets d'encapsulation GRE et expédie un paquet de 1500 octets.
4. Le paquet de 1500 octets ne peut pas traverser la liaison de 1400 octets, il est donc supprimé par le routeur intermédiaire.
5. Le routeur intermédiaire envoie un ICMP (code = 3, type = 4) au routeur GRE avec un MTU de saut suivant de 1400. Le routeur GRE fait passer cette valeur à 1376 (1400 - 24) et définit une valeur interne de MTU IP MTU sur l'interface GRE. Cette modification peut seulement être vue en utilisant la **commande debug tunnel** ; Elle ne peut pas être observée dans la sortie de la commande **show ip interface tunnel<#>**.
6. La fois suivante, l'hôte enverra de nouveau le paquet de 1476 octets, le routeur supprimera le paquet car il est plus volumineux que le MTU IP actuel (1376) de l'interface de tunnel GRE.
7. Le routeur GRE enverra un autre ICMP (code = 3, type = 4) à l'expéditeur avec un MTU de saut suivant de 1376 et l'hôte mettra à jour ses informations à l'aide de la nouvelle valeur.
8. L'hôte envoie de nouveau les données, mais cette fois-ci dans un paquet plus petit de 1376 octets, et GRE ajoute 24 octets d'encapsulation, puis l'achemine. Le paquet est maintenant acheminé vers le partenaire de tunnel GRE, où l'encapsulation est supprimée ; ensuite, il est envoyé à l'hôte de destination.

**Remarque:** Si la commande **tunnel path-mtu-discovery** n'était pas configurée sur le routeur transporteur dans ce scénario, et si le bit DF était défini dans les paquets envoyés via le tunnel GRE, l'hôte 1 parviendrait malgré tout à envoyer les paquets TCP/IP à l'hôte 2, mais en les fragmentant au milieu, au niveau de la liaison MTU 1400. Le partenaire de tunnel GRE devrait également les réassembler avant de pouvoir supprimer l'encapsulation et les acheminer.

## "Mode Tunnel IPsec « pur »

Le protocole de sécurité IP (IPsec) est une méthode basée sur les normes qui apporte confidentialité, intégrité et authenticité aux informations transférées via des réseaux IP; IPsec fournit le cryptage de la couche réseau IP. IPsec rallonge le paquet IP en ajoutant au moins une en-tête IP (mode tunnel). L'en-tête ajoutée varie dans la longueur dépendant le mode de configuration d'IPsec mais elles ne dépassent pas ~58 octets (Protocole ESP (Encapsulating Security Payload) et authentification de l'ESP (ESPauth)) par paquet.

IPsec a deux modes, le mode tunnel et le mode transport.

- Le mode tunnel est le mode par défaut. Avec le mode tunnel, le paquet IP initial entier est protégé (crypté, authentifié, ou les deux) et encapsulé par les en-têtes et les queues de bande IPsec. Une nouvelle en-tête IP est alors ajoutée au début au paquet, spécifiant les points d'extrémité IPsec (homologues) en tant que source et destination. Le mode tunnel peut être utilisé avec tout trafic IP unicast et doit être utilisé si IPsec protège le trafic des hôtes derrière les homologues IPsec. Par exemple, le mode tunnel est utilisé avec des réseaux virtuels privés (VPN) (vpn) sur lesquels les hôtes d'un réseau protégé envoient des paquets aux hôtes d'un autre réseau protégé, via une paire d'homologues IPsec. Avec les VPN, le tunnel IPsec protège le trafic IP entre les hôtes, via le cryptage de ce trafic entre les routeurs homologues IPsec.

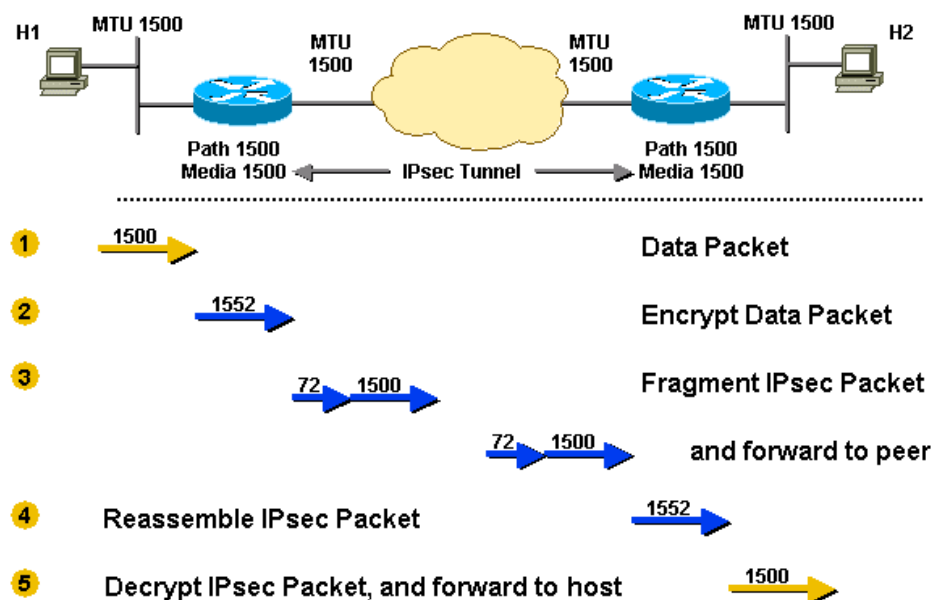
- Avec le mode transport (configuré avec la commande secondaire **mode transport** dans la définition de transformation), seule la charge utile du paquet IP initial est protégée (cryptée, authentifiée, ou les deux). La charge utile est encapsulée par les en-têtes et les remorques d'IPsec. Les en-têtes initiaux d'IP restent intacts, à l'exception de la zone de protocole IP qui est modifiée en ESP (50), et la valeur de protocole d'origine est enregistrée dans la queue de bande IPsec en vue d'une restauration une fois le paquet décrypté. Le mode transport est utilisé seulement quand le trafic IP à protéger est entre les homologues IPsec eux-mêmes, et que les adresses IP source et de destination sont les mêmes que les adresses d'homologues IPsec. Normalement, le mode transport IPsec transport mode est seulement utilisé quand un autre protocole de canalisation en tunnel (comme GRE) est utilisé d'abord pour encapsuler le paquet de données IP, puis IPsec est utilisé pour protéger les paquets de tunnel GRE.

IPsec utilise toujours PMTUD pour les paquets de données et pour ses propres paquets. Certaines commandes de configuration IPsec permettent de modifier le traitement PMTUD pour le paquet IPsec IP, IPsec peut effacer, définir ou copier le bit DF à partir de l'en-tête IP de paquet de données dans l'en-tête IP d'IPsec. Ceci s'appelle « la fonction de remplacement de bit DF ».

**Remarque:** Vous souhaitez éviter la fragmentation après l'encapsulation quand vous faites un cryptage matériel avec IPsec. Le cryptage matériel peut permettre d'obtenir un débit d'environ 50 Mbs, en fonction du matériel, mais si le paquet IPsec est fragmenté, vous perdez de 50 à 90 pour cent du débit. Cette perte s'explique par le fait que les paquets IPsec fragmentés sont commutés pour le réassemblage, puis transmis au moteur de chiffrement, à des fins de décryptement. Cette perte de débit peut mener le débit de cryptage matériel au niveau des performances de cryptage logiciel (2-10 Mbs).

### Scénario 7

Ce scénario décrit la fragmentation IPsec. Dans ce scénario, le MTU du chemin entier est 1500. Dans ce scénario, le bit DF n'est pas défini.

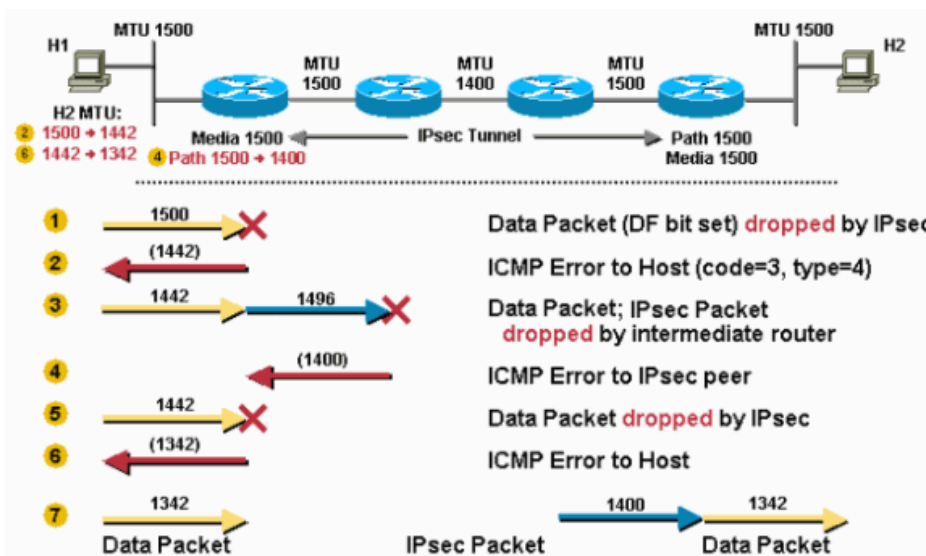


1. Le routeur reçoit un paquet de 1500 octets (en-tête IP de 20 octets + charge utile TCP de 1480 octets) destiné à l'hôte 2.
2. Le paquet de 1500 octets est crypté par IPsec et 52 octets de temps système sont ajoutés (en-tête d'IPsec, queue de bande, et en-tête IP supplémentaire). Maintenant IPsec doit envoyer un paquet de 1552 octets. Puisque le MTU sortant est de 1500, ce paquet devra être fragmenté.
3. Deux fragments sont créés hors du paquet IPsec. Pendant la fragmentation, un en-tête IP de 20 octets supplémentaire est ajouté pour le deuxième fragment, ce qui aboutit à un fragment de 1500 octets et un fragment IP de 72 fragments.
4. Le routeur homologue de tunnel IPsec reçoit les fragments, décolle l'en-tête IP supplémentaire et fusionne les fragments IP dans le paquet IPsec initial. Ensuite, IPsec déchiffre ce paquet.
5. Le routeur achemine ensuite le paquet de données de 1500 octets vers l'hôte 2.

### Scénario 8

Ce scénario est semblable au scénario 6, sauf que dans ce cas, le bit DF est défini dans le paquet de données d'origine ; de plus, une liaison existe sur le chemin entre les homologues de tunnel IPsec portant un MTU plus faible que celui des autres liaisons. Ce scénario explique comment le routeur partenaire d'IPsec exécute les deux rôles PMTUD, comme décrit dans la section Le routeur en tant que participant PMTUD au point d'extrémité d'un tunnel .

Vous verrez dans ce scénario comment PMTU IPsec prend une valeur plus faible en raison du besoin de fragmentation. Rappelez-vous que le bit DF est copié de l'en-tête IP intérieur vers l'en-tête IP extérieur, lorsqu'IPsec crypte un paquet. Les valeurs MTU et PMTU des supports sont stockées dans l'association de sécurité IPsec (SA). Le MTU du support prend comme base le MTU de l'interface du routeur sortant, tandis que le PMTU prend comme base la valeur MTU minimale constatée sur le chemin reliant les homologues IPsec. Rappelez-vous qu'IPsec encapsule/encrypte le paquet avant de tenter de le fragmenter.



1. Le routeur reçoit un paquet de 1500 octets et le supprime car le temps système ISsec, une fois ajouté, rendra le paquet plus volumineux que la valeur de PMTU (1500).
2. Le routeur envoie un message ICMP à l'hôte 1, en indiquant que le MTU du prochain saut est 1442 ( $1500 - 58 = 1442$ ). Ces 58 octets correspondent au temps système IPsec maximum lors de l'utilisation d'IPsec ESP et d'ESPauth. Le temps système IPsec réel peut atteindre 7 octets de moins que cette valeur. L'hôte 1 enregistre cette information, habituellement sous forme de route hôte pour la destination (hôte 2), dans sa table de routage.
3. L'hôte 1 ramène son PMTU pour l'hôte 2 à 1442 ; par conséquent, l'hôte 1 envoie des paquets plus petits (1442 octets) lorsqu'il retransmet les données à l'hôte 2. Le routeur reçoit le paquet de 1442 octets et IPsec ajoute 52 octets de temps système de chiffrement ; le paquet IPsec ainsi obtenu est égal à 1496 octets. Puisque le bit DF de ce paquet est défini dans son en-tête, il est supprimé par le routeur du milieu, avec la liaison MTU 1400 octets.
4. Le routeur du milieu qui a supprimé le paquet envoie un message ICMP à l'expéditeur du paquet IPsec (le premier routeur), indiquant que le MTU du prochain saut est égal à 1400 octets. Cette valeur est enregistrée dans le PMTU IPsec SA.
5. La prochaine fois que l'hôte 1 retransmettra un paquet de 1442 octets (il n'a pas reçu l'accusé de réception correspondant), l'IPsec supprimera le paquet. De nouveau, le routeur supprimera le paquet, car le temps système IPsec, quand il est ajouté au paquet, rendra le paquet plus volumineux que le PMTU (1400).
6. Le routeur envoie un message ICMP à l'hôte 1, en indiquant que le MTU du prochain saut est désormais 1342. ( $1400 - 58 = 1342$ ). L'hôte 1 enregistrera de nouveau ces informations.
7. Quand le hôte 1 retransmet de nouveau les données, il utilise le paquet de plus petite taille (1342). Ce paquet n'exigera pas la fragmentation et l'effectuera via le tunnel IPsec à destination de l'hôte 2.

## GRE et IPsec ensemble

Des interactions plus complexes pour la fragmentation et le PMTUD se produisent quand IPsec est utilisé pour le cryptage des tunnels GRE. IPsec et GRE sont combinés de cette manière parce qu'IPsec ne supporte pas les paquets multicast IP, ce qui signifie que vous ne pouvez pas exécuter un protocole de routage dynamique sur le réseau du VPN d'IPsec. Les tunnels GRE supportent le multicast ; par conséquent, un tunnel GRE peut être utilisé d'abord pour encapsuler le paquet multicast du protocole de routage dans un paquet unicast IP GRE, qui peut ensuite être crypté par IPsec. Dans ce cas, IPsec est souvent déployé en mode transport en plus de GRE, car les homologues IPsec et les points d'extrémité du tunnel GRE (les routeurs) sont identiques, et le mode transport permettra de gagner 20 octets de temps système IPsec.

Un cas intéressant se présente lorsqu'un paquet IP a été divisé en deux fragments et encapsulé par GRE. Dans ce cas, IPsec verra deux GRE indépendants + paquets IP. Souvent, dans les configurations par défaut, l'un de ces paquets est suffisamment volumineux pour nécessiter une fragmentation une fois le cryptage effectué. L'homologue IPsec devra rassembler ce paquet avant déchiffrement. Cette « double fragmentation » (une fois avant GRE et de nouveau après IPsec) sur le routeur expéditeur augmente la latence et diminue le débit. En outre, le réassemblage est commuté par processus ; par conséquent, on assistera à un hit de CPU sur le routeur récepteur chaque fois que cela se produit.

Cette situation peut être évitée en plaçant « ip mtu » sur l'interface de tunnel GRE assez bas pour prendre en considération le temps système de GRE et d'IPsec (par défaut l'interface de tunnel GRE « ip mtu » est définie sur le MTU de l'interface sortante réelle - octets de temps système GRE).

Le tableau suivant présente les valeurs de MTU suggérées pour chaque combinaison tunnel/mode, en supposant que l'interface physique sortante a un MTU de 1500.

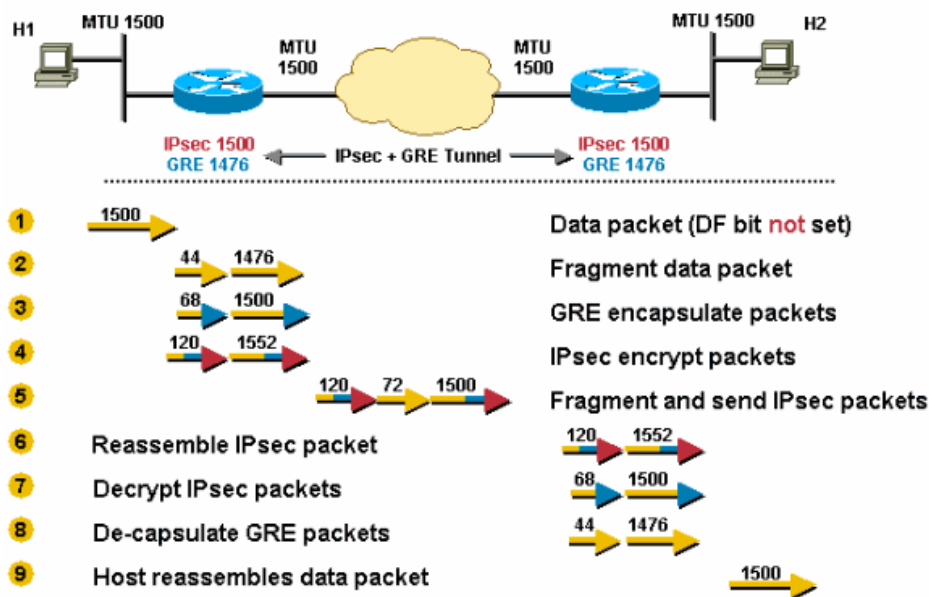
Combinaison de tunnel	MTU spécifique requis	MTU recommandé
GRE + IPsec (transport mode)	1440 bytes	1400 bytes
GRE + IPsec (tunnel mode)	1420 bytes	1400 bytes

**Remarque:** La valeur de MTU de 1400 est recommandée parce qu'elle couvre le plus commun des combinaisons GRE + d'IPsec mode. En outre, il n'y a de du côté incliné pas discernable à tenir compte des 20 ou 40 octets supplémentaires au-dessus. Il est plus facile de se rappeler et de définir une valeur et cette valeur couvre presque tous les scénarios.

## Scénario 9

IPsec est déployé sur GRE. Le MTU physique sortant est 1500, l'IPsec PMTU est 1500, et le GRE est IP MTU 1476 ( $1500 - 24 = 1476$ ). Pour cette raison, les paquets TCP/IP seront fragmentés deux fois, une fois avant GRE et une fois après IPsec. Le paquet sera réduit en fragments avant GRE encapsulation et un de ces GRE packets sera réduit en fragments de nouveau après IPsec encryption.

Configurant « ip mtu 1440 » (IPsec transport mode) ou « ip mtu 1420 » (tunnel IPsec mode) sur le tunnel GRE supprime la possibilité de double fragmentation dans ce scénario.



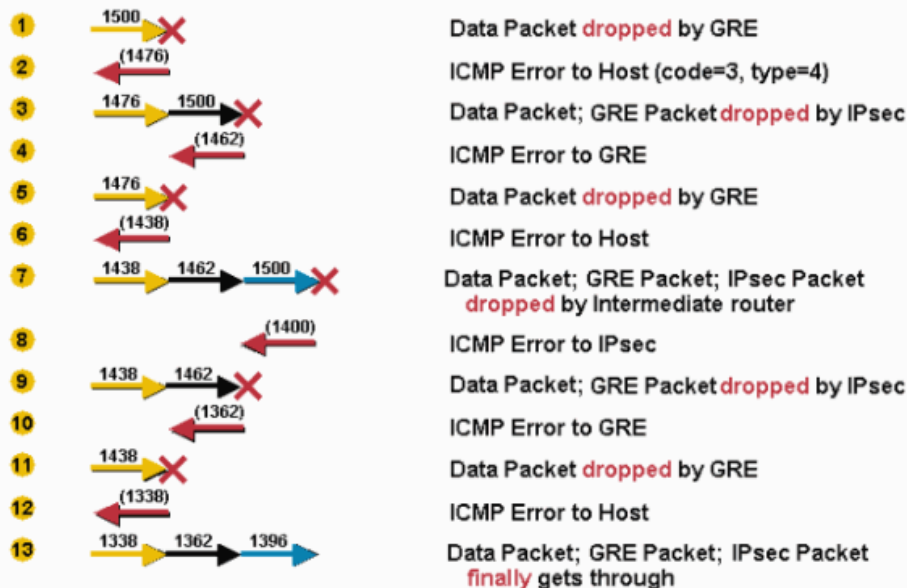
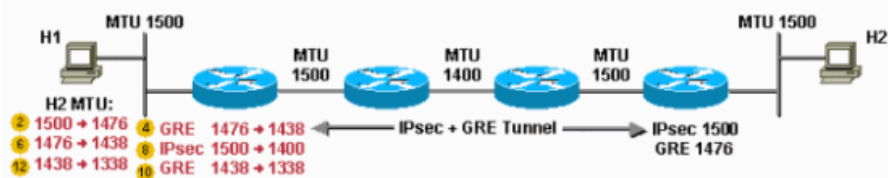
1. Le routeur reçoit un datagramme de 1500 octets.
2. Avant encapsulation, GRE réduit le paquet de 1500 octets dans deux 1476 ( $1500 - 24 = 1476$ ) et 44 (24 données + 20 en-têtes IP) octets de parties.
3. GRE encapsule les fragments d'IP, qui ajoute 24 octets à chaque paquet. Ceci a comme conséquence deux GRE + paquets IPsecs de 1500 ( $1476 + 24 = 1500$ ) et 68 ( $44 + 24$ ) octets chacun.
4. IPsec encrypte les deux paquets, ajoutant 52 octets (mode tunnel d'IPsec) d'encapsulation à chacun, pour obtenir un paquet de 1552 octets et un paquet de 120 octets.
5. Le paquet IPsec de 1552 octets est fragmenté par le routeur, car il est plus grand que le MTU sortant (1500). Le paquet de 1552 octets est divisé en plusieurs parties, un paquet de 1500 octets et un paquet de 72 octets (52 octets de « charge utile » plus une en-tête IP de 20 octets supplémentaires pour le deuxième fragment). Les trois paquets de 1500 octets, 72 octets et 120 octets sont expédiés à IPsec + à l'homologue GRE.
6. Le routeur récepteur réassemble les deux fragments d'IPsec (1500 octets et 72 octets) pour obtenir le paquet d'origine (1552 octets IPsec + GRE). Rien ne doit être fait au paquet de 120 octets IPsec + GRE.
7. IPsec déchiffre les paquets 1552 octets et 120 octets IPsec + GRE afin d'obtenir des paquets de 1500 octets et 68 octets GRE.
8. GRE décapsule les paquets GRE 1500 et 68 octets pour obtenir des fragments de paquets 1476 octets et 44 octets. Ces fragments de paquet IP sont expédiés à l'hôte de destination.
9. L'hôte 2 rassemble ces fragments d'IP pour obtenir le datagramme IP d'origine de 1500 octets.

Le scénario 10 est semblable au scénario 8, à l'exception qu'on note une liaison MTU inférieure dans le chemin de tunnel. C'est un scénario du « pire des cas » pour le premier paquet envoyé de l'hôte 1 à l'hôte 2. Après que la dernière étape dans ce scénario, l'hôte 1 place le PMTU correct pour l'hôte 2 et tout est bien pour les connexions TCP entre l'hôte 1 et les écoulements de TCP de l'hôte 2. entre l'hôte 1 et autre héberge (accessible par l'intermédiaire du tunnel d'IPsec + GRE) devront seulement passer par les trois dernières étapes du scénario 10.

Dans ce scénario, la commande **tunnel path-mtu-discovery** est configurée sur le tunnel GRE et le bit DF est défini sur les paquets TCP/IP en provenance de l'hôte 1.

## Scénario 10





1. Le routeur reçoit un paquet de 1500 octets. Ce paquet est supprimé par GRE, parce que GRE ne peut pas fragmenter ou expédier le paquet car le bit DF est défini et la taille du paquet dépasse l'interface de sortie « ip mtu » après l'ajout du temps système GRE (24 octets).
2. Le routeur envoie un message ICMP à l'hôte 1, en indiquant que le MTU du saut suivant est 1476 (1500 - 24 = 1476).
3. L'hôte 1 modifie son PMTU pour l'hôte 2 en 1476 et envoie la plus petite taille lorsqu'il retransmet le paquet. GRE l'encapsule et remet le paquet de 1500 octets à IPsec. IPsec supprime le paquet parce que GRE a copié le bit DF (défini) à partir de l'en-tête IP interne, et avec le temps système IPsec (maximum 38 octets), le paquet est trop grand pour expédier l'interface physique.
4. IPsec envoie un message ICMP à GRE indiquant que le MTU du saut suivant est égal à 1462 octets (puisqu'un maximum 38 octets sera ajouté pour le cryptage et le temps système IP). GRE enregistre la valeur 1438 (1462 - 24) en tant que « ip mtu » sur l'interface de tunnel.

**Remarque:** Cette variation de la valeur est enregistrée en interne et ne peut pas être vue dans la sortie de la commande **tunnel<#> de show ip interface** . Vous verrez seulement cette modification si vous utilisez la commande **debug tunnel** .

5. La prochaine fois que l'hôte 1 retransmettra le paquet de 1476 octets, GRE le supprimera.
6. Le routeur envoie un message ICMP à l'hôte 1, en indiquant que le MTU du prochain saut est 1438.
7. L'hôte 1 diminue le PMTU pour l'hôte 2 et retransmet un paquet de 1438 octets. Cette fois, GRE accepte le paquet, l'encapsule, et le remet à IPsec pour cryptage. Le paquet IPsec est expédié au routeur intermédiaire et supprimé parce qu'il a un MTU d'interface de sortie de 1400.
8. Le routeur intermédiaire envoie un message ICMP à IPsec lui indiquant que le MTU du prochain saut est égal à 1400. Cette valeur est enregistrée par IPsec dans la valeur PTMU du SA IPsec associé.
9. Quand l'hôte 1 retransmet le paquet de 1438 octets, GRE l'encapsule et le remet à IPsec. IPsec supprime le paquet parce qu'il a changé son propre PMTU en 1400.
10. IPsec envoie un message d'erreur ICMP à GRE en indiquant que le MTU du saut suivant est 1362, et GRE enregistre la valeur 1338 en interne.
11. Quand l'hôte 1 retransmet le paquet d'origine (car il n'a pas reçu d'accusé de réception), GRE le supprime.
12. Le routeur envoie un message ICMP à l'hôte 1, en indiquant que le MTU du prochain saut est 1338 (1500 - 58 - 24 octets). L'hôte 1 abaisse son PMTU pour l'hôte 2 à 1338.
13. L'hôte 1 retransmet un paquet de 1338 octets, et cette fois, il peut finalement accéder complètement à l'hôte 2.

## Autres recommandations

La configuration de la commande **tunnel path-mtu-discovery** sur une interface de tunnel peut favoriser l'interaction entre GRE et IPsec quand elles sont configurées sur le même routeur. Rappelez-vous que si la commande **tunnel path-mtu-discovery** n'est pas configurée, le bit DF est systématiquement effacé dans l'en-tête IP GRE. Cela permet au paquet IP GRE d'être fragmenté, même si l'en-tête IP encapsulé a le bit DF défini, ce qui ne permet pas normalement la fragmentation du paquet.

Si la commande **tunnel path-mtu-discovery** est configurée sur l'interface de tunnel GRE, ce qui suit se produira.

1. GRE copiera le bit DF de l'en-tête IP vers l'en-tête IP GRE.
2. Si le bit DF est défini dans l'en-tête IP GRE et que le paquet est trop volumineux après le cryptage IPsec pour le MTU IP de l'interface physique de sortie, IPsec supprime le paquet et notifie le tunnel GRE en vue d'une réduction de sa taille MTU IP.
3. IPsec applique PMTUD pour ses propres paquets, et si IPsec PMTU change (s'il est réduit), IPsec n'informe pas immédiatement GRE, mais quand un autre paquet « trop grand » se présente, le processus décrit à l'étape 2 se produit.

4. Le MTU IP de GRE est maintenant plus petit ; par conséquent, il supprimera n'importe quel paquet de données IP avec le bit DF défini (devenu trop grand) et enverra un message ICMP à l'hôte expéditeur.

La commande **tunnel path-mtu-discovery** aide commande l'interface GRE à définir dynamiquement son MTU IP, plutôt que de le faire de façon statique avec la commande **ip mtu** . Il est recommandé d'utiliser ces deux commandes. La commande **ip mtu** est utilisée pour fournir au GRE l'espace et le temps système d'IPsec relatifs au MTU IP de l'interface physique sortante locale. La commande **tunnel path-mtu-discovery** permet de diminuer davantage le MTU IP du tunnel GRE s'il existe une liaison MTU IP dans le chemin entre les homologues IPsec.

Voici ce que vous pouvez faire si vous rencontrez des problèmes PMTUD au sein d'un réseau sur lequel GRE + des tunnels IPsec sont configurés.

La liste suivante commence par la solution la plus souhaitable.

- Réglez le problème de non fonctionnement de PMTUD, qui est habituellement provoqué par le blocage d'ICMP par un routeur ou un pare-feu.
- Utilisez la commande **ip tcp adjust-mss** sur les interfaces de tunnel afin que le routeur diminue la valeur MSS TCP au niveau du paquet TCP SYN. Cela aide les deux hôtes d'extrémité (l'expéditeur et le destinataire TCP) à utiliser des paquets suffisamment petits pour ne pas nécessiter l'intervention de PMTUD.
- Utilisez le routage spécifique sur l'interface d'entrée du routeur et configurez une route pour effacer le bit DF de l'en-tête IP avant l'accès à l'interface de tunnel GRE. Ceci permettra la fragmentation du paquet de données IP avant l'encapsulation GRE.
- Augmentez la valeur « ip mtu » sur l'interface de tunnel GRE pour qu'elle soit égale à la valeur MTU de l'interface de sortie. Ceci permettra au paquet de données IP d'être encapsulé sans fragmentation préalable. Le paquet GRE sera alors crypté par IPsec, puis fragmenté pour la sortie de l'interface physique de sortie. Dans ce cas, vous ne configurez pas la commande **tunnel path-mtu-discovery** sur l'interface de tunnel GRE. Ceci peut diminuer considérablement le débit, car le réassemblage de paquets IP sur l'homologue IPsec s'effectue en mode de commutation de processus.

## Informations connexes

- [Découverte de MTU RFC 1191](#)
- [Options de découverte RFC 1063 IP MTU](#)
- [Internet Protocol RFC 791](#)
- [RFC 793 Transmission Control Protocol \(TCP\)](#)
- [RFC 879 Taille maximale des segments TCP et rubriques connexes](#)
- [RFC 1701 Encapsulation de routage générique \(GRE\)](#)
- [RFC 1241 Plan pour protocole IEP \(Internet Encapsulation Protocol\)](#)
- [Encapsulation RFC 2003 IP dans IP](#)
- [Livre blanc Technologies](#)

---

© 1992-2010 Cisco Systems Inc. Tous droits réservés.

---

Date du fichier PDF généré: 19 octobre 2014

---

[http://www.cisco.com/cisco/web/support/CA/fr/109/1094/1094907\\_pmtud\\_ipfrag.html](http://www.cisco.com/cisco/web/support/CA/fr/109/1094/1094907_pmtud_ipfrag.html)

---